

# UNCERTAINTY QUANTIFICATION FOR REGRESSION USING PROPER SCORING RULES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Quantifying uncertainty of machine learning model predictions is essential for reliable decision-making, especially in safety-critical applications. Recently, uncertainty quantification (UQ) theory has advanced significantly, building on a firm basis of learning with proper scoring rules. However, these advances were focused on classification, while extending these ideas to regression remains challenging. In this work, we introduce a unified UQ framework for regression based on proper scoring rules, such as CRPS, logarithmic, squared error, and quadratic scores. We derive closed-form expressions for the resulting uncertainty measures under practical parametric assumptions and show how to estimate them using ensembles of models. In particular, the derived uncertainty measures naturally decompose into aleatoric and epistemic components. The framework recovers popular regression UQ measures based on predictive variance and differential entropy. Our broad evaluation on synthetic and real-world regression datasets provides guidance for selecting reliable UQ measures.

## 1 INTRODUCTION

Predictive models trained with machine learning are widely used, yet their predictions are often uncertain. Quantifying this uncertainty is essential, especially in safety-critical applications (Helou et al., 2020; Su et al., 2023). Uncertainty quantification (UQ) provides a way to measure the confidence of a model in its predictions (Hüllermeier & Waegeman, 2021; Gal & Ghahramani, 2016). One of the central goals of UQ is to separate aleatoric uncertainty (AU) associated with irreducible noise inherent to the data-generating process, from epistemic uncertainty (EU), which arises from the limited knowledge of the underlying model (Hüllermeier & Waegeman, 2021). Distinguishing these types can help identify the causes of uncertainty and select the mitigation strategy.

Although UQ of machine learning models predictions is a mature field, it was historically mostly relying on ad hoc heuristic uncertainty measures. Recently, the coherent theoretical body was developed (Wimmer et al., 2023; Kotelevskii et al., 2025; Hofman et al., 2024; Schweighofer et al., 2025), mainly focusing on classification tasks. The core idea of these methods is to consider pointwise risk as a natural measure of predictive uncertainty. Early works (Kotelevskii et al., 2022; Lahlou et al., 2023) applied this idea to specific losses and model families. More recent works extended this approach via proper scoring rules (Gneiting & Raftery, 2007), deriving general uncertainty measures through Bayesian approximations to the different components of the pointwise risk (Kotelevskii et al., 2025; Hofman et al., 2024; Schweighofer et al., 2025).

However, many real-world problems involve the prediction of continuous, unbounded outcomes (Gal & Ghahramani, 2016; Amini et al., 2020; Valdenegro-Toro & Mori, 2022), while aforementioned approaches have not yet been adapted for regression. This paper closes that gap by extending the proper-scoring-rule UQ framework to regression tasks. We consider various appropriate scoring rules developed for regression, such as CRPS which is a leading tool in meteorology Alet et al. (2025); Lang et al. (2025). That allows us to derive a wide family of uncertainty measures that can be readily decomposed into aleatoric and epistemic components. Interestingly, our framework allows to recover some of the established entropy- and variance-based measures for UQ in regression (Depeweg et al., 2018; Bülte et al., 2025).

Our main **contributions** can be summarized as follows.

- A theoretical formulation of regression UQ based on proper scoring rules, with separate quantification of aleatoric and epistemic components; see Section 3
- Closed-form, practically computable approximations for these components under standard regression assumptions for ensemble-based estimators; see Tables 2,3 and Appendix A.
- A broad empirical evaluation on synthetic and real-world datasets, covering selective prediction, out-of-distribution detection, and active learning, culminating in recommendations for selecting uncertainty measures in practice; see Section 5.

## 2 BACKGROUND ON PROPER SCORING RULES

We start by defining the proper scoring rules and introducing the necessary notation. This section is based on the related works on proper scoring rules (Gneiting & Raftery, 2007; Waghmare & Ziegel, 2025), and follows their notation.

We consider supervised learning problems with input  $x \in \mathcal{X}$  and output  $y \in \mathcal{Y}$ . Let  $\mathcal{P}$  be a convex class of probability measures on  $\mathcal{Y}$ . A probabilistic forecast is any probability measure  $\hat{P} \in \mathcal{P}$ . For classification, machine learning models usually produce a probabilistic forecast by default: the output is the distribution  $\hat{p} = \hat{p}(y | x)$  over class labels. In regression problems, the standard approach is to proceed with just a point prediction  $\hat{y} = f(x) \in \mathcal{Y}$ . However, to efficiently deal with predictive uncertainty, a model may output parameters of the distribution  $\hat{p}(y | x)$ , e.g.,  $\hat{\mu}(x), \hat{\sigma}(x)$  for a Gaussian distribution.

A scoring rule is an extended real-valued function  $S: \mathcal{P} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}} = [-\infty, +\infty]$  such that  $S(\hat{P}, \cdot)$  is integrable for all  $\hat{P} \in \mathcal{P}$ . Thus if the forecast is  $\hat{P}$  and  $y$  materializes, the forecaster’s penalty is  $S(\hat{P}, y)$ . Usually, we are interested in assessing the predictive quality of the forecast  $\hat{P}$  on average over possible values of  $y$ . To achieve that, we consider the expected score under  $P$ :

$$S(\hat{P}, P) = \int S(\hat{P}, y) dP(y) \quad (1)$$

The practically important family of scoring rules are so-called *proper scoring rules* that incentivize the learning of the true data distribution. The scoring rule  $S$  is called proper relative to  $\mathcal{P}$  if

$$S(P, P) \leq S(\hat{P}, P) \text{ for all } \hat{P}, P \in \mathcal{P}. \quad (2)$$

The scoring rule is strictly proper if the equality in (2) holds iff  $\hat{P} = P$ . Thus, any (strictly) proper scoring can be used for evaluation of probabilistic forecasts ensuring that its minimum is attained at the true data distribution.

The following definitions regarding proper scoring rules are useful for our future exposition.

**Definition 1.** An expected score between a measure and itself is called the *entropy function*:

$$H(P) = S(P, P) = \int S(P, y) dP(y).$$

Entropy determines the smallest possible error that one can achieve which corresponds to the Bayes optimal predictor (the true data distribution).

**Definition 2.** A *divergence function* is the difference between the expected scores of the predicted and true distributions:

$$d(\hat{P}, P) = S(\hat{P}, P) - S(P, P) = S(\hat{P}, P) - H(P),$$

quantifying how much worse it is to predict with distribution  $\hat{P}$  than with the true data distribution  $P$ .

**Proper scoring rules for regression.** A number of proper scoring rules exist for probability distributions with continuous support (Gneiting & Raftery, 2007). In this work we consider several of the most common proper scoring rules that could be used for regression such as *continuous ranked probability score (CRPS)*, *logarithmic*, *quadratic and squared error scores* (see their definitions in Table 1). All of these scores are strictly proper except for the squared error score, which is just proper. In what follows, we will show how to derive practical uncertainty estimates for regression based on these scores.

Table 1: Examples of commonly considered proper scoring rules for regression.

Name	Definition
Continuous ranked probability score	$\text{CRPS}(\hat{P}, y) = \int_{\mathbb{R}} (F_{\hat{P}}(t) - \mathbb{I}\{y \leq t\})^2 dt$
Logarithmic score	$\text{LS}(\hat{P}, y) = -\log \hat{p}(y)$
Quadratic score	$\text{QS}(\hat{P}, y) = -2\hat{p}(y) + \int_{\mathbb{R}} \hat{p}(t)^2 dt$
Squared error score	$\text{SE}(\hat{P}, y) = (y - \mathbb{E}_{Y \sim \hat{P}}[Y])^2$

### 3 REGRESSION UNCERTAINTY MEASURES VIA PROPER SCORING RULES

In this section, we introduce the uncertainty quantification framework for regression and show how one can end up with practical equations for uncertainty measures using proper scoring rules.

#### 3.1 POINTWISE RISK AS AN UNCERTAINTY MEASURE

Let  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  be a training set of  $n$  independent and identically distributed (i.i.d.) samples from the joint distribution  $P(X, Y)$ , where  $X_i \in \mathcal{X} = \mathbb{R}^d$  and  $Y_i \in \mathcal{Y} = \mathbb{R}$ . We denote the true conditional distribution of  $Y$  at  $X = x$  by  $P(Y | X = x)$ . The standard goal of statistical learning is to approximate this conditional distribution with some model  $\hat{P}(Y | X = x)$  using the training data. For brevity, we further assume that everything is implicitly conditional on  $X = x$ , so we write  $P = P(Y | X = x)$  and  $\hat{P} = \hat{P}(Y | X = x)$ . We will denote the corresponding densities by  $p$  and  $\hat{p}$ , and the corresponding cumulative distribution functions by  $F_P$  and  $F_{\hat{P}}$ .

The alternative view on the expected score  $S(\hat{P}, P)$  in (1) is to interpret it as an expected error of prediction by model  $\hat{P}$  at a given point. Such an object is usually called *pointwise risk* in statistical learning literature. There are three key risk quantities that we consider in our framework:

- **Total risk.** The expected score between the estimated  $\hat{P}$  and the true  $P$  distributions, capturing the overall uncertainty of a prediction:

$$R_{\text{Tot}}(\hat{P}, P) = \int S(\hat{P}, y) dP(y) = S(\hat{P}, P).$$

- **Bayes risk.** The risk we would incur if we could predict with the *true* distribution itself:

$$R_{\text{Bayes}}(P) = \int S(P, y) dP(y) = S(P, P) = H(P).$$

Since it does not depend on the model, the nature of this error is purely *aleatoric*. It is equal to the entropy of the data distribution and, for any proper score, it equals to the smallest possible expected error of prediction.

- **Excess risk.** It measures how much worse the model is compared to the perfect prediction:

$$R_{\text{Exc}}(\hat{P}, P) = R_{\text{Tot}}(\hat{P}, P) - R_{\text{Bayes}}(P) = d(\hat{P}, P).$$

It is equal to the score’s divergence and captures the part of the error not described by the randomness in the data itself. Thus, it is naturally related to *epistemic* uncertainty.

The risks above depend on the true distribution  $P$ , which is unknown. In the next section, we will consider various approximations of these risks.

#### 3.2 APPROXIMATIONS OF THE RISKS

Following Kotelevskii et al. (2025); Hofman et al. (2024); Schweighofer et al. (2025), we consider Bayesian approximations to the risk expressions. We assume that we have a predictive parametric model, with the vector of parameters  $\theta$  following a posterior distribution of  $p(\theta | \mathcal{D})$ , giving a *distribution* over predictive distributions  $\hat{P}_{\theta}$ .

While the framework is model-agnostic, in order to obtain tractable closed-form expressions for particular risk approximations we proceed to make practical parametric assumptions about the

Table 2: Bayes-risk approximations under different proper scoring rules. For the Gaussian mixture model, we use the following notation:  $\hat{P}_{\text{ens}} = \frac{1}{M} \sum_{i=1}^M \mathcal{N}(\mu_i, \sigma_i^2)$ . For  $\hat{R}_{\text{Bayes}}^3$ , different plug-in estimates of the variance  $\sigma_*^2$  can be used (see Appendix A.1.2).

App.	CRPS	Logarithmic	Quadratic	SE
$\hat{R}_{\text{Bayes}}^1$	$\frac{1}{M\sqrt{\pi}} \sum_{i=1}^M \sigma_i$	$\frac{1}{2M} \sum_{i=1}^M \log(2\pi e\sigma_i^2)$	$-\frac{1}{2M\sqrt{\pi}} \sum_{i=1}^M \frac{1}{\sigma_i}$	$\frac{1}{M} \sum_{i=1}^M \sigma_i^2$
$\hat{R}_{\text{Bayes}}^2$	$\frac{1}{2M^2} \sum_{i=1}^M \sum_{j=1}^M A(\mu_{ij}, \sigma_{ij})$	$-\int_{\mathbb{R}} \hat{p}_{\text{ens}}(y) \log \hat{p}_{\text{ens}}(y) dy$	$-\int_{\mathbb{R}} \hat{p}_{\text{ens}}(y)^2 dy$	$\frac{1}{M} \sum_{i=1}^M (\sigma_i^2 + \mu_i^2) - \mu_*^2$
$\hat{R}_{\text{Bayes}}^3$	$\frac{1}{\sqrt{\sigma_*^2/\pi}}$	$\frac{1}{2} \log(2\pi e\sigma_*^2)$	$-\frac{1}{2\sqrt{\pi}\sigma_*}$	$\sigma_*^2$

Table 3: Excess risks  $\hat{R}_{\text{Exc}}^{1,1}$  (Bayesian averaging of both  $P$  and  $\hat{P}$ ) for different proper scoring rules.

Scoring Rule	$\hat{R}_{\text{Exc}}^{1,1}$
CRPS	$\frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \left[ A(\mu_{ij}, \sigma_{ij}) - \frac{\sigma_i + \sigma_j}{\sqrt{\pi}} \right]$
Logarithmic	$\frac{1}{2M^2} \sum_{i=1}^M \sum_{j=1}^M \left[ \frac{\sigma_j^2 + (\mu_i - \mu_j)^2}{\sigma_i^2} - 1 \right]$
Quadratic	$-\frac{2}{M} \sum_{i=1}^M H(\hat{P}_i) - \frac{2}{M^2} \sum_{i=1}^M \sum_{j=1}^M \mathcal{N}(\mu_i   \mu_j, \sigma_i^2 + \sigma_j^2)$
SE	$2\widehat{\text{Var}}[\mu_i] = \frac{2}{M} \sum_{i=1}^M (\mu_i - \mu_*)^2$

distribution  $\hat{P}_\theta$ . For the scope of this work, we consider the Gaussian assumption  $\hat{P}_\theta = \mathcal{N}(\mu, \sigma^2)$ , i.e. that the predictive distribution is defined by a vector of distribution parameters  $\theta = (\mu, \sigma^2)$ . Other parametric assumptions, e.g. Laplace, could be used to derive different measures. As a Bayesian model  $p(\theta | \mathcal{D})$  we consider an ensemble of Gaussians with parameters  $\theta_i = (\mu_i, \sigma_i^2)$ :  $\hat{P}_{\text{ens}} = \frac{1}{M} \sum_{i=1}^M \hat{P}_{\theta_i} = \frac{1}{M} \sum_{i=1}^M \mathcal{N}(\mu_i, \sigma_i^2)$ .

From this perspective, we estimate any of the risks above in the following generic ways:

- Bayesian averaging of the risk.** Posterior expectation of the risk:  $\mathbb{E}_{p(\theta|\mathcal{D})} [R(\hat{P}, \hat{P}_\theta)]$ .
- Posterior predictive distribution.** First, average the distributions  $\hat{P}_{\text{ens}} = \mathbb{E}_{p(\theta|\mathcal{D})} [\hat{P}_\theta]$ , and use the plug-in as an estimate of the risk:  $R(\hat{P}, \hat{P}_{\text{ens}})$ .
- Gaussian surrogate.** We consider two Gaussian approximations for the posterior  $\hat{P}_*, \bar{P}_*$  (see Appendix A.1.2), which are used as plug-ins for the risk estimation:  $R(\hat{P}, \hat{P}_* \text{ or } \bar{P}_*)$ .

Note that either of these approaches can also be used to derive a specific  $\hat{P}$  estimate as well (Kotelevskii et al., 2025; Hofman et al., 2024; Schweighofer et al., 2025).

**Practical approximation.** Mixing the ‘‘how to treat  $P$ ’’ choices with the ‘‘what we feed in as  $\hat{P}$ ’’ choices yields (in general) nine different concrete approximation pairs  $(\hat{P}, P)$ . We discuss what particular options we choose in A.2 and summarize it in Table 7. Therefore, any risk estimate considered from now on carries two superscripts  $\hat{R}^{i,j}$ , stating these choices. Numbers correspond to the generic way to approximate risks we outlined above, e.g., 1 denotes Bayesian averaging of risks. For Bayes risk, we need to build only the approximation of the true distribution  $P$ , for which we have three options. We illustrate the results for popular regression proper scoring rules for the approximation of Bayes risk in Table 2, and  $\hat{R}_{\text{Exc}}^{1,1}$  in Table 3. We refer to Appendix A for complete derivations and other approximation options<sup>1</sup>.

**Generalization over earlier work.** Bülte et al. (2025) recently examined entropy- and variance-based uncertainty measures for regression under a new set of axioms. Those uncertainty measures are

<sup>1</sup>We use the following notation:  $\mathcal{N}(a | b, s^2) = \frac{1}{\sqrt{2\pi s^2}} e^{-\frac{(a-b)^2}{2s^2}}$ ,  $A(\mu, \sigma) = 2\sigma\phi\left(\frac{\mu}{\sigma}\right) + \mu [2\Phi\left(\frac{\mu}{\sigma}\right) - 1]$ , and  $\Phi, \phi$  are CDF and PDF of a standard normal.

216 defined as (Depeweg et al., 2018):  
 217

$$218 \quad \mathbb{H}(\mathbb{E}_{p(\theta|\mathcal{D})}[\hat{P}_\theta]) = \mathbb{E}_{p(\theta|\mathcal{D})}[\mathbb{H}(\hat{P}_\theta)] + \mathbb{E}_{p(\theta|\mathcal{D})}[d_{\text{KL}}(\hat{P}_\theta, \mathbb{E}_{p(\theta|\mathcal{D})}[\hat{P}_\theta])], \quad (3)$$

$$219 \quad \text{Var}(\mathbb{E}_{p(\theta|\mathcal{D})}[\hat{P}_\theta]) = \mathbb{E}_{p(\theta|\mathcal{D})}[\text{Var}(\hat{P}_\theta)] + \text{Var}_{p(\theta|\mathcal{D})}[\mathbb{E}(\hat{P}_\theta)] \quad (4)$$

220  
 221  
 222 respectively, where  $\mathbb{H}$  is the Shannon entropy and  $d_{\text{KL}}$  is the Kullback–Leibler divergence. In our  
 223 scoring-rule view, these measures arise naturally as the *logarithmic* and *SE* special cases (using  
 224  $\hat{R}_{\text{Tot}}^{1,2} = \hat{R}_{\text{Bayes}}^1 + \hat{R}_{\text{Exc}}^{1,2}$  as risk approximations), see the derivations in Appendix A.2.2 and A.2.4 and  
 225 Example 2.1 in Bülte et al. (2025). Thus, our proper scoring rule framework thus comprises entropy-  
 226 and variance-based measures and introduces additional measures based on CRPS, a score widely  
 227 used in meteorology for probabilistic forecasts, and the quadratic score. Furthermore, it allows us to  
 228 easily derive new measures under different parametric assumptions or using alternative scoring rules.  
 229

## 230 4 RELATED WORK

231  
 232  
 233 **Axiomatic approaches in classification.** Despite the maturity of uncertainty quantification, the  
 234 formalization of uncertainty measures appeared only recently. Most of the efforts in this formalization  
 235 were made in the context of classification, where authors tried to formalize what makes a “good”  
 236 uncertainty score. Specifically, Wimmer et al. (2023); Sale et al. (2024) propose different axiom sets  
 237 for classification, but there seems to be no existing measure that satisfies all of them.  
 238

239 **Risk-based uncertainty measures.** The idea of viewing predictive uncertainty as *pointwise risk*,  
 240 the expected value of a loss (statistical pointwise risk), was first put forward for classification  
 241 in (Kotelevskii et al., 2022) and both settings in (Lahlou et al., 2023). However, they considered a  
 242 specific class of models and specific loss functions. Building on this view, in (Schweighofer et al.,  
 243 2025; Kotelevskii et al., 2025; Hofman et al., 2024) showed that the risk of any proper scoring rule  
 244 admits a clean decomposition into the *Bayes risk* (aleatoric part) and the *excess risk* (epistemic part),  
 245 where the former corresponds to the generalized entropy, while the latter is the notion of a Bregman  
 246 divergence (Bregman, 1967). However, generalizing the ideas of proper scoring rules to regression  
 247 was not considered.

248 **Towards regression.** The recent paper Bülte et al. (2025) introduces axioms for regression and  
 249 assesses entropy and variance for them. However, they do not introduce other measures that can be  
 250 derived from considering proper scoring rules. Our paper fills this gap by deriving regression-specific  
 251 scores from proper scoring rules, specifically, the CRPS score.

252 **Practical deep-learning baselines.** Bayesian dropout (Gal & Ghahramani, 2016), Deep Ensem-  
 253 bles (Lakshminarayanan et al., 2017), and the heteroscedastic losses (Kendall & Gal, 2017) remain  
 254 standard tools of uncertainty quantification in regression in Deep Learning, because they scale  
 255 well, though they lack a firm axiomatic footing. Deep Evidential Regression (Amini et al., 2020)  
 256 offers a single-pass alternative that outputs Normal-Inverse-Gamma parameters, giving closed-form  
 257 aleatoric and epistemic variances, yet recent studies question how reliably DER captures epistemic  
 258 risk (Juergens et al., 2024).  
 259

## 260 5 EXPERIMENTS

261  
 262  
 263 In our experiments, we investigate the empirical behavior of the uncertainty measures introduced  
 264 in Section 3. We begin by examining how these measures respond to different perturbations of the  
 265 posterior  $p(\theta | \mathcal{D})$  and showcase their behavior on a synthetic regression task. We then evaluate  
 266 their utility on three downstream problems: selective prediction, out-of-distribution detection, and  
 267 active learning. To further assess the relationship between measures, we evaluate rank correlations  
 268 between measures and examine when different scores lead to different decisions. Unless stated  
 269 otherwise, uncertainties are computed by Monte Carlo approximation through heteroscedastic deep  
 ensembles (Lakshminarayanan et al., 2017); training details are provided in Appendix E.1.

### 5.1 CHARACTERIZATION OF UNCERTAINTY MEASURES

We investigate the behavior of the different uncertainty measures by applying shifts to the posterior  $p(\theta | \mathcal{D})$ , where  $\theta = (\mu, \sigma^2)$ , inspired by the investigation in Bülte et al. (2025, Figure 2). The original posterior is a uniform distribution in the predicted means  $\mu \sim \mathcal{U}(-1, 1)$  and a uniform distribution in the predicted variances  $\sigma^2 \sim \mathcal{U}(1, 2)$ . This closed-form posterior allows us to estimate uncertainty measures directly without relying on sampling methods such as deep ensembles. We consider four different shifts to the posterior: (a) a location shift on the predicted means, (b) a location shift on the predicted variances, (c) a scale shift of the predicted means, and (d) a scale shift of the predicted variances, see Figure 1 (left column).

**Location shift of predicted means.** The distribution of predicted means shifts to  $\mu \sim \mathcal{U}(1, 3)$ , thus the location of the posterior shifts, but not its scale. The results are shown in Figure 1 (first row). None of the measures considered changes under this shift. This is desirable, as uncertainty measures should not be sensitive to the magnitude of mean predictions.

**Location shift of predicted variances.** The distribution of predicted variances shifts to  $\sigma^2 \sim \mathcal{U}(2, 3)$ , thus the location of the posterior shifts, but not its scale. The results are shown in Figure 1 (second

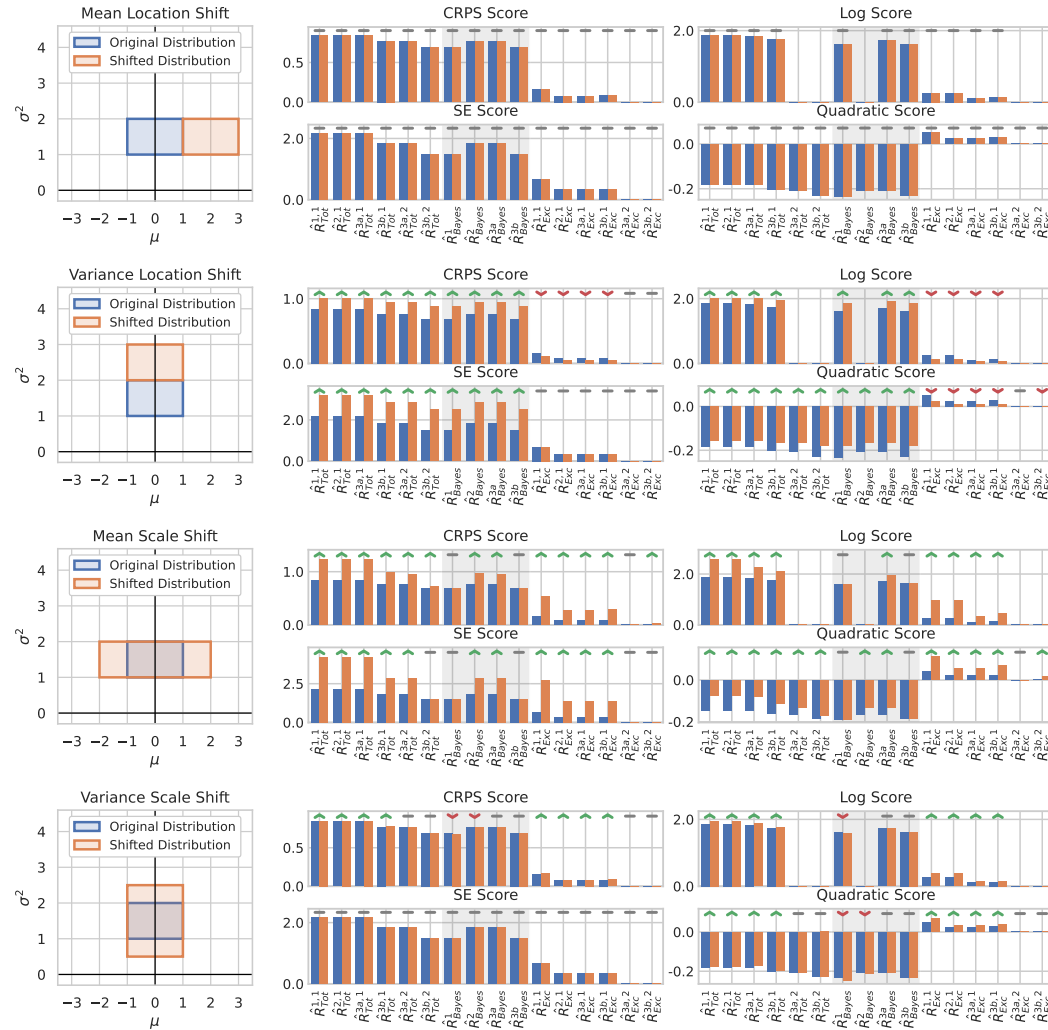


Figure 1: Behavior of uncertainty measures under location and scale shifts of the posterior distribution. Arrows indicate whether a measure increased or decreased due to the shift, gray bars indicate changes  $< 1\%$ , missing entries that the measure is not computable or constant zero.

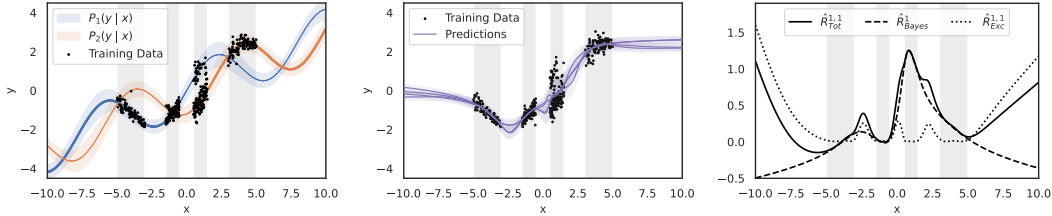


Figure 2: Uncertainty with the logarithmic scoring rule for the synthetic task. **Left:** The ground truth generative distribution is a mixture of two distributions. The thickness of the line indicates the weight of the mixture component. Intervals where training data is sampled are highlighted in gray. **Middle:** Predictive distributions of different models trained on the training data. **Right:** Resulting Total, Bayes and Excess risks.

row). It is expected that all uncertainty measures increase under such a shift. However, we observe that excess risk measures either decrease or stay the same (for  $\hat{R}_{\text{Exc}}^{3a,2}$  under CRPS and Quadratic and for all excess risks under SE score). Here, staying the same is a desirable outcome over decreasing.

**Scale shift of predicted means.** The distribution of predicted means shifts to  $\mu \sim \mathcal{U}(-2, 2)$ , thus the scale of the posterior shifts, but not its location. The results are shown in Figure 1 (third row). It is expected that all uncertainty measures increase under such a shift. We observe that  $\hat{R}_{\text{Bayes}}^1$  and  $\hat{R}_{\text{Bayes}}^{3b}$  under all scores are invariant under such a shift, the same for  $\hat{R}_{\text{Tot}}^{3b,2}$  under SE. This could reduce their viability in downstream tasks. All other measures increase under the scale shift in the predicted means as expected.

**Scale shift of predicted variances.** The distribution of predicted variances shifts to  $\sigma^2 \sim \mathcal{U}(0.5, 2.5)$ , thus the scale of the posterior shifts, but not its location. The results are shown in Figure 1 (fourth row). It is expected that total and epistemic uncertainty increase under such a shift, but that the aleatoric uncertainty stays the same.

**General observations.** Across all shifts, we observe that  $\hat{R}_{\text{Exc}}^{3a,2}$  and  $\hat{R}_{\text{Exc}}^{3b,2}$  have very low magnitudes, and for the SE score, they are equal to zero. Furthermore, the magnitude of the excess risks is often much lower than the Bayes risks, and they scale most with the mean scale shift, which is expected.

## 5.2 SYNTHETIC DATASETS

Having characterized how the measures react to controlled shifts of the posterior, we now turn to a synthetic data example that verifies these behaviors against common intuition. Specifically, inputs far from the training distribution should have greater epistemic uncertainty, while regions where the data-generating process is inherently ambiguous should yield higher aleatoric uncertainty.

In our 1D regression experiment (see results in Figure 2 and its description in Appendix E.2), the *excess risk* estimate, a proxy for epistemic uncertainty, increases for inputs outside the training data support, where the model cannot reliably infer the underlying dependency. The *Bayes risk* estimate, capturing aleatoric uncertainty, is largest where the mapping from  $x$  to  $y$  is least determined. In our synthetic example, this occurs in the region where the probabilities of sampling from the two generating curves are similar. Consequently, the *total risk* is high whenever either component is high.

For clarity, we present a Log-score instantiation (which leads to entropy-based measures) in the main text and visualize the corresponding approximations of Bayes, excess, and total risk in Figure 2. Log-score-based uncertainty measures are widely used in practice (Bülte et al., 2025), and our empirical results support this common choice. We also observe the same qualitative behavior for other instantiations, with additional visualizations provided in Appendix E.2.

## 5.3 SELECTIVE PREDICTION

In selective prediction, the model’s performance is evaluated on a specific subset that is defined as having low uncertainty. Therefore, the uncertainty measures the ability to indicate whether or not the prediction is correct. We considered multiple datasets for these experiments; see Appendix E.3

Table 4: Results for selective prediction (PRR  $\downarrow$ ) for different scoring rules (SR), averaged over datasets. The best result is underlined, all within its standard deviation in bold.

SR	$\hat{R}_{Tot}^{1,1}$	$\hat{R}_{Tot}^{2,1}$	$\hat{R}_{Tot}^{3a,1}$	$\hat{R}_{Tot}^{3b,1}$	$\hat{R}_{Tot}^{3a,2}$	$\hat{R}_{Tot}^{3b,2}$	$\hat{R}_{Bayes}^1$	$\hat{R}_{Bayes}^2$	$\hat{R}_{Bayes}^{3a}$	$\hat{R}_{Bayes}^{3b}$	$\hat{R}_{Exc}^{1,1}$	$\hat{R}_{Exc}^{2,1}$	$\hat{R}_{Exc}^{3a,1}$	$\hat{R}_{Exc}^{3b,1}$	$\hat{R}_{Exc}^{3a,2}$	$\hat{R}_{Exc}^{3b,2}$
CRPS	<b>0.318</b> $\pm 0.007$	<b>0.318</b> $\pm 0.007$	<b>0.319</b> $\pm 0.007$	0.327 $\pm 0.007$	0.327 $\pm 0.007$	0.352 $\pm 0.007$	0.356 $\pm 0.007$	0.327 $\pm 0.007$	0.327 $\pm 0.007$	0.357 $\pm 0.007$	0.339 $\pm 0.006$	0.339 $\pm 0.006$	0.340 $\pm 0.006$	0.341 $\pm 0.006$	0.504 $\pm 0.006$	0.378 $\pm 0.006$
Log	<b>0.323</b> $\pm 0.007$	<b>0.323</b> $\pm 0.007$	<b>0.321</b> $\pm 0.007$	0.327 $\pm 0.007$	-	-	0.356 $\pm 0.007$	-	0.327 $\pm 0.007$	0.357 $\pm 0.007$	0.397 $\pm 0.006$	0.397 $\pm 0.006$	0.400 $\pm 0.006$	0.398 $\pm 0.006$	-	-
SE	<b>0.320</b> $\pm 0.007$	<b>0.320</b> $\pm 0.007$	<b>0.320</b> $\pm 0.007$	0.327 $\pm 0.007$	0.327 $\pm 0.007$	0.357 $\pm 0.007$	0.357 $\pm 0.007$	0.327 $\pm 0.007$	0.327 $\pm 0.007$	0.357 $\pm 0.007$	<b>0.325</b> $\pm 0.006$	<b>0.325</b> $\pm 0.006$	<b>0.325</b> $\pm 0.006$	<b>0.325</b> $\pm 0.006$	-	-
Quad.	<b>0.319</b> $\pm 0.007$	<b>0.319</b> $\pm 0.007$	<b>0.323</b> $\pm 0.007$	0.327 $\pm 0.007$	0.329 $\pm 0.007$	0.347 $\pm 0.007$	0.356 $\pm 0.007$	0.326 $\pm 0.007$	0.327 $\pm 0.007$	0.357 $\pm 0.007$	0.514 $\pm 0.007$	0.514 $\pm 0.007$	0.523 $\pm 0.006$	0.511 $\pm 0.007$	0.689 $\pm 0.004$	0.503 $\pm 0.007$

Table 5: Results for out-of-distribution detection (AUROC  $\uparrow$ ) for different scoring rules (SR), averaged over datasets. The best result is underlined, all within its standard deviation in bold.

SR	$\hat{R}_{Tot}^{1,1}$	$\hat{R}_{Tot}^{2,1}$	$\hat{R}_{Tot}^{3a,1}$	$\hat{R}_{Tot}^{3b,1}$	$\hat{R}_{Tot}^{3a,2}$	$\hat{R}_{Tot}^{3b,2}$	$\hat{R}_{Bayes}^1$	$\hat{R}_{Bayes}^2$	$\hat{R}_{Bayes}^{3a}$	$\hat{R}_{Bayes}^{3b}$	$\hat{R}_{Exc}^{1,1}$	$\hat{R}_{Exc}^{2,1}$	$\hat{R}_{Exc}^{3a,1}$	$\hat{R}_{Exc}^{3b,1}$	$\hat{R}_{Exc}^{3a,2}$	$\hat{R}_{Exc}^{3b,2}$
CRPS	0.794 $\pm 0.012$	0.794 $\pm 0.012$	0.794 $\pm 0.012$	0.779 $\pm 0.013$	0.777 $\pm 0.012$	0.727 $\pm 0.009$	0.714 $\pm 0.010$	0.777 $\pm 0.012$	0.777 $\pm 0.012$	0.714 $\pm 0.010$	<b>0.825</b> $\pm 0.011$	<b>0.825</b> $\pm 0.011$	<b>0.825</b> $\pm 0.011$	<b>0.825</b> $\pm 0.011$	0.795 $\pm 0.006$	<b>0.826</b> $\pm 0.010$
Log	0.802 $\pm 0.011$	0.802 $\pm 0.011$	0.797 $\pm 0.012$	0.785 $\pm 0.013$	-	-	0.713 $\pm 0.010$	-	0.777 $\pm 0.012$	0.714 $\pm 0.010$	<b>0.827</b> $\pm 0.011$	<b>0.827</b> $\pm 0.011$	<b>0.826</b> $\pm 0.011$	<b>0.826</b> $\pm 0.011$	-	-
SE	0.792 $\pm 0.012$	0.792 $\pm 0.012$	0.792 $\pm 0.012$	0.777 $\pm 0.012$	0.777 $\pm 0.012$	0.714 $\pm 0.010$	0.714 $\pm 0.010$	0.777 $\pm 0.012$	0.777 $\pm 0.012$	0.714 $\pm 0.010$	<b>0.820</b> $\pm 0.010$	<b>0.820</b> $\pm 0.010$	<b>0.820</b> $\pm 0.010$	<b>0.820</b> $\pm 0.010$	-	-
Quad.	0.800 $\pm 0.012$	0.800 $\pm 0.012$	0.800 $\pm 0.012$	0.784 $\pm 0.013$	0.777 $\pm 0.013$	0.740 $\pm 0.011$	0.713 $\pm 0.010$	0.777 $\pm 0.012$	0.777 $\pm 0.012$	0.714 $\pm 0.010$	<b>0.816</b> $\pm 0.016$	<b>0.816</b> $\pm 0.016$	<b>0.816</b> $\pm 0.016$	<b>0.817</b> $\pm 0.015$	0.783 $\pm 0.005$	<b>0.822</b> $\pm 0.012$

Table 6: Results for active learning (average rank  $\downarrow$  over datasets and five seeds using different scores as acquisition functions). The best result is underlined, all within its standard deviation in bold.

SR	$\hat{R}_{Tot}^{1,1}$	$\hat{R}_{Tot}^{2,1}$	$\hat{R}_{Tot}^{3a,1}$	$\hat{R}_{Tot}^{3b,1}$	$\hat{R}_{Tot}^{3a,2}$	$\hat{R}_{Tot}^{3b,2}$	$\hat{R}_{Bayes}^1$	$\hat{R}_{Bayes}^2$	$\hat{R}_{Bayes}^{3a}$	$\hat{R}_{Bayes}^{3b}$	$\hat{R}_{Exc}^{1,1}$	$\hat{R}_{Exc}^{2,1}$	$\hat{R}_{Exc}^{3a,1}$	$\hat{R}_{Exc}^{3b,1}$	$\hat{R}_{Exc}^{3a,2}$	$\hat{R}_{Exc}^{3b,2}$	Random
CRPS	16.60 $\pm 3.57$	16.60 $\pm 3.57$	16.57 $\pm 4.31$	17.83 $\pm 3.26$	17.60 $\pm 4.19$	18.27 $\pm 2.75$	19.07 $\pm 4.38$	16.60 $\pm 2.75$	17.70 $\pm 3.64$	19.20 $\pm 3.04$	17.20 $\pm 3.58$	17.20 $\pm 3.58$	15.90 $\pm 3.68$	15.07 $\pm 4.01$	15.27 $\pm 3.00$	15.87 $\pm 3.53$	20.50 $\pm 5.89$
Log	15.73 $\pm 2.63$	15.73 $\pm 2.63$	14.97 $\pm 2.58$	15.47 $\pm 3.60$	-	-	19.13 $\pm 2.20$	-	17.97 $\pm 3.61$	21.47 $\pm 2.55$	<b>13.90</b> $\pm 1.53$	<b>13.90</b> $\pm 1.33$	17.33 $\pm 3.97$	<b>14.37</b> $\pm 3.93$	-	-	20.50 $\pm 5.89$
SE	16.90 $\pm 4.80$	16.90 $\pm 4.80$	16.90 $\pm 4.80$	16.87 $\pm 3.26$	16.67 $\pm 3.03$	16.47 $\pm 3.91$	16.47 $\pm 3.91$	16.67 $\pm 3.03$	16.67 $\pm 3.03$	16.47 $\pm 3.91$	17.50 $\pm 3.67$	17.50 $\pm 3.67$	16.00 $\pm 3.21$	16.00 $\pm 3.21$	-	-	20.50 $\pm 5.89$
Quad.	29.23 $\pm 2.06$	29.23 $\pm 2.06$	32.57 $\pm 2.34$	29.47 $\pm 5.12$	31.70 $\pm 4.93$	31.27 $\pm 4.23$	31.07 $\pm 3.07$	32.47 $\pm 4.10$	30.33 $\pm 3.81$	31.83 $\pm 4.17$	16.10 $\pm 3.37$	16.10 $\pm 3.37$	<b>12.60</b> $\pm 1.97$	<b>13.43</b> $\pm 3.93$	17.00 $\pm 2.34$	15.90 $\pm 2.46$	20.50 $\pm 5.89$

for details. Empirical findings for selective prediction in classification (e.g., Kotelevskii et al., 2025; Schweighofer et al., 2025) suggest that total uncertainty performs well in this task, which we seek to validate for the regression setting as well.

The results are shown in Table 4, which are averages over all the considered datasets (see Table 9 for detailed results). Performances are measured as prediction-reject ratios (PRRs), where a lower PRR indicates more effective rejection of inaccurate predictions using the uncertainty score (details in Appendix E.3). The best performing scores are  $\hat{R}_{Tot}^{1,1}$ ,  $\hat{R}_{Tot}^{2,1}$  and  $\hat{R}_{Tot}^{3a,1}$  for all considered scoring rules. Furthermore, all excess risks under the SE score, which are equivalent to each other, perform similarly well as the total risks. Noteworthy, excess risks for the quadratic score are performing very poorly compared to all other measures; the same for  $\hat{R}_{Exc}^{3a,2}$  for the CRPS scoring rule. Furthermore, excess risks under the widely considered Log-score perform **slightly** worse than most other measures. In general, we recommend using  $\hat{R}_{Tot}^{1,1}$  for selective prediction. It leads to the overall best performance under CRPS, and is consistently among the best across scoring rules.

#### 5.4 OUT-OF-DISTRIBUTION DETECTION

Out-of-distribution (OOD) detection is widely considered as task for uncertainty estimation. We want to detect OOD inputs, as there is no guarantee that a model trained on some in-distribution (ID) dataset will perform well on it. An uncertainty score often indicates that a new input is OOD, assigning high uncertainty to such inputs. We considered a dataset consisting of a mosaic of four MNIST (Lecun et al., 1998) images, where the target is given as the number formed by those four digits, as an ID dataset. We considered multiple OOD datasets; for details, see Appendix E.4.

This task is evaluated using the AUROC of the uncertainty score to distinguish between ID and OOD data. Results are provided in Table 5, averaged over all considered ID and OOD data pairs. We refer to Table 10 in the Appendix for detailed results per ID/OOD pair. We find that excess risk measures are generally the most suitable for this task, particularly  $\hat{R}_{Exc}^{1,1}$ . This aligns with the assumptions on excess risk in prior work (Lahlou et al., 2023; Hofman et al., 2024; Kotelevskii et al., 2025).

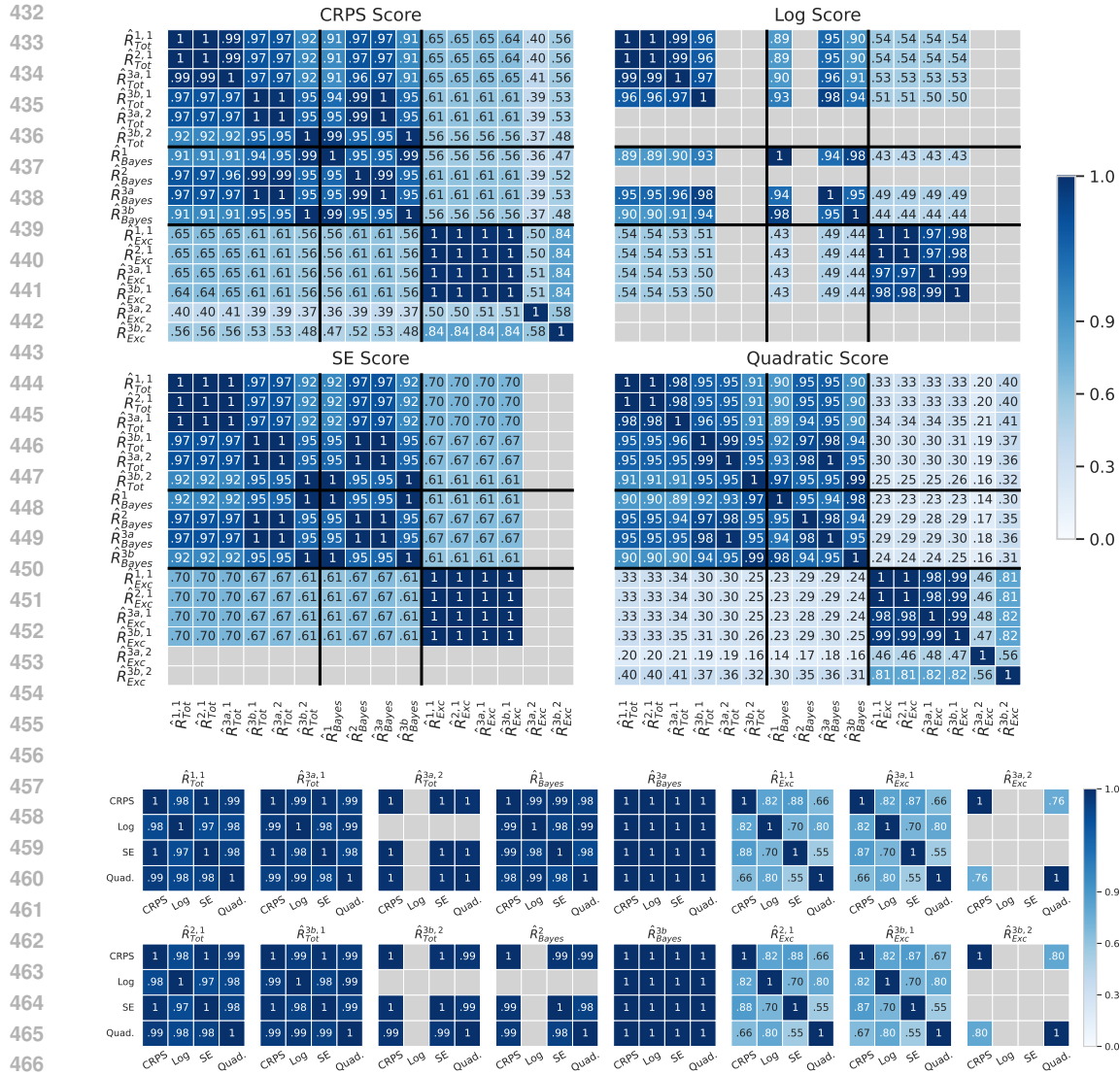


Figure 3: Kendall’s  $\tau_b$  rank correlation between different risk approximations for different scoring rules (top row) and between different scoring rules for the considered risk approximations (bottom rows). Correlations are averaged over all considered datasets.

### 5.5 ACTIVE LEARNING

Another widely considered real-world task utilizing uncertainty information is active learning. The measure of uncertainty should be used to guide which samples are most uncertain, and thus interesting to obtain labels to add to the training dataset. Generally, excess risks (epistemic uncertainty) are considered to be better guidance than total or Bayes risks (Mukhoti et al., 2023).

The results are provided in Table 6, showing average ranks over the considered datasets; see Appendix E.5 for details. As expected, we observe that excess risks are the best performing measures, but there is no clear best approximation strategy. Most outperform random sampling, yet there are some exceptions: the total and Bayes risks under the quadratic score and  $\hat{R}_{Bayes}^{3b}$  under the log score.

### 5.6 RANK CORRELATION BETWEEN MEASURES

Finally, we inspect the rank correlation between the measures within our framework. We use the same datasets and models as for selective prediction and calculate the rank correlations induced by

different uncertainty measures. The results shown in Figure 3 consider the rank correlation between different risk approximations per scoring rule, as well as the rank correlations between scoring rules given a risk approximation. The most important findings are, that  $\widehat{R}_{\text{Bayes}}^1$  is very strongly correlated to  $\widehat{R}_{\text{Bayes}}^{3b}$  and the same for  $\widehat{R}_{\text{Bayes}}^2$  and  $\widehat{R}_{\text{Bayes}}^{3a}$ . Furthermore,  $\widehat{R}_{\text{Tot}}^{1,1}$  and  $\widehat{R}_{\text{Tot}}^{2,1}$  are equivalent due to the linearity in the expectation, the same holds for  $\widehat{R}_{\text{Exc}}^{1,1}$  and  $\widehat{R}_{\text{Exc}}^{2,1}$ . Regarding differences between induced rankings for different scoring rules for the same approximations, a striking finding is that for both  $\widehat{R}_{\text{Bayes}}^{3a}$  and  $\widehat{R}_{\text{Bayes}}^{3b}$ , all scoring rules are perfectly correlated, the same for  $\widehat{R}_{\text{Tot}}^{3a,2}$  for those scoring rules it is computable for. Furthermore, rank correlations are overall very high among approximations of total and Bayes risks and only substantially vary across excess risks. **Although total and Excess risk estimates for the scores are not perfectly uncorrelated, performance on the previous tasks clearly shows the usefulness of utilizing different risks for specific tasks. A recent benchmark of Mucsányi et al. (2024) provides a perspective on this in classification scenarios.**

## 6 DISCUSSION

**General findings and failure modes.** Throughout our considered experiments, we found that the particular choice of proper scoring rule and approximation strategy do not have a major impact in that they would lead to much worse results. The exception is  $\bar{R}_{\text{Exc}}^{3a,2}$ , which we found performs significantly worse throughout all considered experiments. Furthermore, this particular risk approximation was found to have undesirable characteristics across scoring rules, such as invariance under mean scale shift, linking theoretical behavior to downstream performance. We hypothesize, that this is caused by the fact that  $\bar{R}_{\text{Exc}}^{3a,2}$  uses two approximation strategies where the disagreement between individual models is averaged out. Additionally, we found that no measure of Bayes risk, regardless of scoring rule or approximation strategy, is best in any of the considered tasks.

**Recommendations.** Our empirical findings in Section 5 show that different risks are desirable as measures in different tasks, leading to the following recommendations. For selective prediction, we find that total risk leads to the best results across particular scoring rules and approximation strategies. For both out-of-distribution detection and active learning, we find that excess risks perform best. While most scoring rules and Bayesian approximation strategies give reasonable results, we recommend  $\bar{R}_{\text{Tot}}^{1,1}$  and  $\bar{R}_{\text{Exc}}^{1,1}$  for any of the considered scoring rules as a default choice.

## 7 CONCLUSION

We presented a unified framework for uncertainty quantification in regression based on proper scoring rules. By extending the risk-based decomposition of uncertainty from classification to regression, we derived principled measures of total, aleatoric, and epistemic uncertainty that generalize and subsume existing variance- and entropy-based approaches. Our framework admits closed-form expressions under common assumptions, provides flexible Bayesian and surrogate approximations, and yields practical estimators suitable for deep ensembles.

Through a broad empirical evaluation, we demonstrated that these measures behave consistently under controlled perturbations, align with theoretical desiderata, and provide competitive or superior performance in uncertainty related tasks. Importantly, we highlighted both commonalities and differences across scoring rules, offering guidance for choosing uncertainty measures in practice.

In this work, we demonstrate our general framework under the Gaussian parametric assumption, using Deep Ensembles to compute MC approximations of posterior expectations to estimate the risks. We chose this particular instantiation due to their widespread use in practice. Looking into the possibility to derive closed-form solutions under different scoring rules for different parametric assumptions and using different approaches for Bayesian approximation, together with investigating their empirical properties, is an important direction for future work.

Overall, our work bridges the gap between theoretical considerations and practical deep learning methods, providing a principled and versatile foundation for regression uncertainty quantification.

## REPRODUCIBILITY STATEMENT

We provide the full code to reproduce our experiments as supplementary material and will release it publicly upon acceptance. Furthermore, we added detailed descriptions of all experiments in Appendix E. All experiments were conducted on publicly available datasets or datasets we created ourselves, which will be released alongside the code. We ran experiments with multiple seeds, if applicable, and report summary statistics.

## USAGE OF LARGE LANGUAGE MODELS (LLMs)

LLMs were used as a general-purpose assistive tool during the preparation of this paper. Their usage fell into two categories: (i) for writing assistance, they helped improve clarity and readability of certain passages through language refinement and (ii) for coding assistance, where they provided support with code completion and debugging. LLMs were not used for research ideation, experimental design, theoretical development, or analysis of results. All substantive contributions, including the conception of ideas, methodology, and experiments, were made by the authors.

## REFERENCES

- Ferran Alet, Ilan Price, Andrew El-Kadi, Dominic Masters, Stratis Markou, Tom R Andersson, Jacklynn Stott, Remi Lam, Matthew Willson, Alvaro Sanchez-Gonzalez, et al. Skillful joint probabilistic weather forecasting from marginals. *arXiv preprint arXiv:2506.10772*, 2025.
- Alexander Amini, Wilko Schwarting, Ava Soleimany, and Daniela Rus. Deep evidential regression. *Advances in neural information processing systems*, 33:14927–14937, 2020.
- Rafael Ballester-Ripoll, Enrique G. Paredes, and Renato Pajarola. Sobol tensor trains for global sensitivity analysis. *Reliability Engineering & System Safety*, 2019.
- Thierry Bertin-Mahieux, Daniel Ellis, Brian Whitman, and Paul Lamere. The million song dataset. *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 2011.
- Lev M Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR computational mathematics and mathematical physics*, 7(3):200–217, 1967.
- Christopher Bülte, Yusuf Sale, Timo Löhr, Paul Hofman, Gitta Kutyniok, and Eyke Hüllermeier. An axiomatic assessment of entropy-and variance-based uncertainty quantification in regression. *arXiv preprint arXiv:2504.18433*, 2025.
- Krisztian Buza. Feedback prediction for blogs. In *Data analysis, machine learning and knowledge discovery*, pp. 145–152. Springer, 2013.
- Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. EMNIST: Extending MNIST to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*. IEEE, 2017.
- Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Finale Doshi-Velez, and Steffen Udluft. Decomposition of Uncertainty in Bayesian Deep Learning for Efficient and Risk-sensitive Learning. In *Proceedings of the 35th International Conference on Machine Learning*. PMLR, 2018.
- Jean-Louis Durrieu, Jean-Philippe Thiran, and Finnian Kelly. Lower and upper bounds for approximation of the kullback-leibler divergence between gaussian mixture models. *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4833–4836, 2012.

- 594 Kelwin Fernandes, Pedro Vinagre, and Paulo Cortez. A proactive intelligent decision support system  
595 for predicting the popularity of online news. In *Progress in Artificial Intelligence*, Cham, 2015.  
596 Springer International Publishing.
- 597 Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model  
598 uncertainty in deep learning. In *International conference on machine learning*, pp. 1050–1059.  
599 PMLR, 2016.
- 600  
601 Tilmann Gneiting. Making and evaluating point forecasts. *Journal of the American Statistical*  
602 *Association*, 106:746 – 762, 2009.
- 603  
604 Tilmann Gneiting and Adrian E. Raftery. Strictly proper scoring rules, prediction, and estimation.  
605 *Journal of the American Statistical Association*, 102:359 – 378, 2007.
- 606  
607 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
608 recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,  
2016.
- 609  
610 Marieka A Helou, Deborah DiazGranados, Michael S Ryan, and John W Cyrus. Uncertainty in  
611 decision making in medicine: a scoping review and thematic analysis of conceptual models.  
612 *Academic Medicine*, 95(1):157–165, 2020.
- 613  
614 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint*  
*arXiv:1606.08415*, 2016.
- 615  
616 John R. Hershey and Peder A. Olsen. Approximating the kullback leibler divergence between  
617 gaussian mixture models. *2007 IEEE International Conference on Acoustics, Speech and Signal*  
*Processing - ICASSP '07*, 4:IV–317–IV–320, 2007.
- 618  
619 Paul Hofman, Yusuf Sale, and Eyke Hüllermeier. Quantifying aleatoric and epistemic uncertainty  
620 with proper scoring rules. *arXiv preprint arXiv:2404.12215*, 2024.
- 621  
622 Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning:  
623 An introduction to concepts and methods. *Machine learning*, 110(3):457–506, 2021.
- 624  
625 Alexander Immer, Emanuele Palumbo, Alexander Marx, and Julia Vogt. Effective bayesian het-  
626 eroscedastic regression with deep neural networks. *Advances in Neural Information Processing*  
*Systems*, 36:53996–54019, 2023.
- 627  
628 Mira Juergens, Nis Meinert, Viktor Bengs, Eyke Hüllermeier, and Willem Waegeman. Is epistemic  
629 uncertainty faithfully represented by evidential deep learning methods? In *International Conference*  
*on Machine Learning*, pp. 22624–22642. PMLR, 2024.
- 630  
631 Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer  
632 vision? *Advances in neural information processing systems*, 30, 2017.
- 633  
634 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
*arXiv:1412.6980*, 2014.
- 635  
636 Andreas Kirsch, Joost van Amersfoort, and Yarin Gal. BatchBALD: Efficient and diverse batch  
637 acquisition for deep bayesian active learning. In *Advances in Neural Information Processing*  
*Systems*, 2019.
- 638  
639 Nikita Kotelevskii, Aleksandr Artemenkov, Kirill Fedyanin, Fedor Noskov, Alexander Fishkov,  
640 Artem Shelmanov, Artem Vazhentsev, Aleksandr Petiushko, and Maxim Panov. Nonparametric  
641 uncertainty quantification for single deterministic neural network. *Advances in Neural Information*  
*Processing Systems*, 35:36308–36323, 2022.
- 642  
643 Nikita Kotelevskii, Vladimir Kondratyev, Martin Takáč, Eric Moulines, and Maxim Panov. From  
644 risk to uncertainty: Generating predictive uncertainty measures via bayesian estimation. In *The*  
645 *Thirteenth International Conference on Learning Representations*, 2025.
- 646  
647 Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.  
*University of Toronto*, 2009.

- 648 Salem Lahlou, Moksh Jain, Hadi Nekoei, Victor I Butoi, Paul Bertin, Jarrid Rector-Brooks, Maksym  
649 Korablyov, and Yoshua Bengio. Deup: Direct epistemic uncertainty prediction. *Transactions on*  
650 *Machine Learning Research*, 2023.
- 651 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
652 uncertainty estimation using deep ensembles. *Advances in neural information processing systems*,  
653 30, 2017.
- 654 Simon Lang, Martin Leutbecher, and Pedro Maciel. A multi-scale loss formulation for learning a  
655 probabilistic model with proper score optimisation. *arXiv preprint arXiv:2506.10868*, 2025.
- 656 Yann Lecun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to  
657 document recognition. *Proceedings of the IEEE*, 1998.
- 658 Andrey Malinin and Mark Gales. Uncertainty estimation in autoregressive structured prediction. In  
659 *International Conference on Learning Representations*, 2021.
- 660 Bálint Mucsányi, Michael Kirchhof, and Seong Joon Oh. Benchmarking uncertainty disentanglement:  
661 Specialized uncertainties for specialized tasks. *Advances in neural information processing systems*,  
662 37:50972–51038, 2024.
- 663 Jishnu Mukhoti, Andreas Kirsch, Joost van Amersfoort, Philip H.S. Torr, and Yarin Gal. Deep  
664 deterministic uncertainty: A new simple baseline. In *Proceedings of the IEEE/CVF Conference on*  
665 *Computer Vision and Pattern Recognition*, 2023.
- 666 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al.  
667 Reading digits in natural images with unsupervised feature learning. *NIPS workshop on deep*  
668 *learning and unsupervised feature learning*, 2011.
- 669 David A. Nix and Andreas S. Weigend. Estimating the mean and variance of the target probability dis-  
670 tribution. In *Proceedings of 1994 IEEE International Conference on Neural Networks (ICNN'94)*,  
671 1994.
- 672 Prashant Rana. Physicochemical properties of protein tertiary structure. UCI Machine Learning  
673 Repository, 2013.
- 674 Yusuf Sale, Viktor Bengs, Michele Caprio, and Eyke Hüllermeier. Second-order uncertainty quan-  
675 tification: A distance-based approach. In *International Conference on Machine Learning*, pp.  
676 43060–43076. PMLR, 2024.
- 677 Leonard J. Savage. Elicitation of personal probabilities and expectations. *Journal of the American*  
678 *Statistical Association*, 66:783–801, 1971.
- 679 Kajetan Schweighofer, Lukas Aichberger, Mykyta Ielanskyi, and Sepp Hochreiter. On information-  
680 theoretic measures of predictive uncertainty. *Uncertainty in Artificial intelligence*, 2025.
- 681 Maximilian Seitzer, Arash Tavakoli, Dimitrije Antic, and Georg Martius. On the pitfalls of het-  
682 eroscedastic uncertainty estimation with probabilistic neural networks. In *International Conference*  
683 *on Learning Representations*, 2022.
- 684 Nicki Skafté, Martin Jørgensen, and Søren Hauberg. Reliable training and estimation of variance  
685 networks. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2019.
- 686 Sanbao Su, Yiming Li, Sihong He, Songyang Han, Chen Feng, Caiwen Ding, and Fei Miao. Un-  
687 certainty quantification of collaborative detection for self-driving. In *2023 IEEE International*  
688 *Conference on Robotics and Automation (ICRA)*, pp. 5588–5594. IEEE, 2023.
- 689 Pinar Tüfekci. Prediction of full load electrical power output of a base load operated combined cycle  
690 power plant using machine learning methods. *International Journal of Electrical Power & Energy*  
691 *Systems*, 2014.
- 692 Matias Valdenegro-Toro and Daniel Saromo Mori. A deeper look into aleatoric and epistemic  
693 uncertainty disentanglement. In *2022 IEEE/CVF Conference on Computer Vision and Pattern*  
694 *Recognition Workshops (CVPRW)*, pp. 1508–1516. IEEE, 2022.

702 Kris Villez. Analytical expressions to compute the continuous ranked probability score (crps).  
703 Technical report, Swiss Federal Institute of Aquatic Science and Technology, 2017.  
704

705 Kartik Waghmare and Johanna Ziegel. Proper scoring rules for estimation and forecast evaluation.  
706 *ArXiv*, 2504.01781, 2025.

707 Lisa Wimmer, Yusuf Sale, Paul Hofman, Bernd Bischl, and Eyke Hüllermeier. Quantifying aleatoric  
708 and epistemic uncertainty in machine learning: Are conditional entropy and mutual information  
709 appropriate measures? In *Uncertainty in artificial intelligence*, pp. 2282–2292. PMLR, 2023.  
710

711 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking  
712 machine learning algorithms. *ArXiv*, 1708.07747, 2017.  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A RISK ESTIMATION FOR PROPER SCORING RULES

This appendix develops the general machinery for estimating pointwise risk under proper scoring rules in regression. It formalizes Bayes/total/excess risks and derives tractable approximation for CRPS, logarithmic, quadratic, and squared error scores under a Gaussian assumption.

### A.1 BAYESIAN RISK ESTIMATION

This section specifies Bayesian estimators for the three risk components, showing how to combine posterior averaging, posterior-predictive (mixture) plug-ins, and Gaussian surrogates to obtain computable expressions.

#### A.1.1 GAUSSIAN ENSEMBLE MODEL

We will consider a specific tractable example where our Bayesian model is a finite Gaussian ensemble. This model arises e.g. when we train multiple neural network copies  $\{\hat{P}_i\}_{i=1}^M$  that predict Gaussian parameters starting from different weights initializations and is sometimes called *deep ensembles* (Lakshminarayanan et al., 2017):

$$\hat{P}_i = \mathcal{N}(\mu_i, \sigma_i^2), \quad \mu_i \in \mathbb{R}, \sigma_i \in \mathbb{R}_+, i = 1, \dots, M.$$

We can interpret this as a finite i.i.d. sample drawn from a posterior distribution over parameters of a one-dimensional Gaussian. Our prediction is the posterior predictive, which equals the full Gaussian mixture distribution:

$$\hat{P}_{\text{ens}} = \frac{1}{M} \sum_{i=1}^M \mathcal{N}(\mu_i, \sigma_i^2), \quad \hat{p}_{\text{ens}}(y) = \frac{1}{M} \sum_{i=1}^M \hat{p}_i(y).$$

#### A.1.2 ENSEMBLE APPROXIMATIONS WITH A SINGLE GAUSSIAN

Previous work on ensembles (Lakshminarayanan et al., 2017) has used the following Gaussian approximation to the full mixture distribution:

$$\hat{P}_* = \mathcal{N}(\mu_*, \sigma_*^2), \quad \mu_* = \frac{1}{M} \sum_{i=1}^M \mu_i, \quad \sigma_*^2 = \frac{1}{M} \sum_{i=1}^M (\sigma_i^2 + \mu_i^2) - \mu_*^2, \quad (5)$$

where  $\mu_*$  and  $\sigma_*^2$  are the mean and variance of the full mixture distribution. Another, more simple Gaussian approximation is to set the variance equal to the average variance of the components (Bülte et al., 2025):

$$\bar{P}_* = \mathcal{N}(\mu_*, \bar{\sigma}_*^2), \quad \mu_* = \frac{1}{M} \sum_{i=1}^M \mu_i, \quad \bar{\sigma}_*^2 = \frac{1}{M} \sum_{i=1}^M \sigma_i^2. \quad (6)$$

### A.2 TRACTABLE RISK ESTIMATION USING GAUSSIAN ENSEMBLES

We assume that we only have access to a single Gaussian ensemble model and propose to use in to approximate both  $\hat{P}$  and  $P$ . Based on the results of the previous section, we arrive at the following four approximations: (1) Bayesian risk averaging, (2) full mixture distribution  $\hat{P}_{\text{ens}}$ , and (3)  $\hat{P}_*$ . A specific risk estimate can thus be denoted as  $\hat{R}^{i,j}$ ,  $i, j = 1, \dots, 3$ . These cases roughly correspond to the ones for classification that appeared in (Kotelevskii et al., 2025). We summarize them in Table 7, where empty cells correspond to combinations that do not seem particularly useful (approximation for ground truth should be at least as complex as for the prediction). In what follows, we will provide expressions for these estimates for different proper scoring rules.

#### A.2.1 CRPS SCORE

The notation and derivations of various quantities related to CRPS for the ensemble model can be found in the Appendix B.1. In particular, we will use the following notation:

$$A(\mu, \sigma) := 2\sigma\phi\left(\frac{\mu}{\sigma}\right) + \mu\left[2\Phi\left(\frac{\mu}{\sigma}\right) - 1\right], \quad \mu_{ij} := \mu_i - \mu_j, \quad \sigma_{ij} := \sqrt{\sigma_i^2 + \sigma_j^2}.$$

Table 7: Chosen combinations of estimates for excess risk estimation.

$\downarrow \hat{P}, \rightarrow P$	BA	$\hat{P}_{\text{ens}}$	$\hat{P}_*$
BA	$\hat{R}^{1,1}$	-	-
$\hat{P}_{\text{ens}}$	$\hat{R}^{2,1}$	-	-
$\hat{P}_*$	$\hat{R}^{3,1}$	$\hat{R}^{3,2}$	-

**Bayes risk.** Since Bayes risk is equal to the entropy of the true distribution  $H(P)$ , we only need to choose an approximation scheme for  $P$ :

- **Bayesian averaging:**  $\hat{R}_{\text{Bayes}}^1 = \frac{1}{M} \sum_{i=1}^M H(\hat{P}_i) = \frac{1}{M\sqrt{\pi}} \sum_{i=1}^M \sigma_i$ .
- **Posterior predictive:**  $\hat{R}_{\text{Bayes}}^2 = H(\hat{P}_{\text{ens}}) = \frac{1}{2} \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M A(\mu_{ij}, \sigma_{ij})$ .
- **Gaussian approximation:**  $\hat{R}_{\text{Bayes}}^3 = H(\hat{P}_*) = \sqrt{\frac{\sigma_*^2}{\pi}}$ .

**Excess risk.** Since the divergence function for CRPS is symmetric, many of  $\hat{R}_{\text{Exc}}^{i,j}$  quantities coincide. Here we present all possible combinations of approximations for the Excess risk in the Gaussian ensemble model.

- **Bayesian averaging for both.** Also known as *expected pairwise Bregman divergence*.

$$\hat{R}_{\text{Exc}}^{1,1} = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M d(\hat{P}_i, \hat{P}_j) = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \left[ A(\mu_{ij}, \sigma_{ij}) - \frac{\sigma_i + \sigma_j}{\sqrt{\pi}} \right] \quad (7)$$

$$= \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M A(\mu_{ij}, \sigma_{ij}) - \frac{2}{M} \sum_{i=1}^M \frac{\sigma_i}{\sqrt{\pi}}. \quad (8)$$

- **Gaussian mixture and Bayesian averaging.** This approximation is also referred to as *Bregman information* (and *reverse Bregman information* for non-symmetric divergences).

$$\hat{R}_{\text{Exc}}^{2,1} = \hat{R}_{\text{Exc}}^{1,2} = \frac{1}{M} \sum_{i=1}^M d(\hat{P}_i, \hat{P}_{\text{ens}}) = \frac{1}{M} \sum_{i=1}^M \left[ \text{CRPS}(\hat{P}_i, \hat{P}_{\text{ens}}) - H(\hat{P}_{\text{ens}}) \right] \quad (9)$$

$$= \frac{1}{2} \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M A(\mu_{ij}, \sigma_{ij}) - \frac{1}{M} \sum_{i=1}^M \frac{\sigma_i}{\sqrt{\pi}}. \quad (10)$$

- **Moment-matched Gaussian approximation and Bayesian averaging:**

$$\hat{R}_{\text{Exc}}^{3,1} = \hat{R}_{\text{Exc}}^{1,3} = \frac{1}{M} \sum_{j=1}^M d(\hat{P}_*, \hat{P}_j) = \frac{1}{M} \sum_{j=1}^M \left[ A(\mu_{*j}, \sigma_{*j}) - \frac{\sigma_* + \sigma_j}{\sqrt{\pi}} \right] \quad (11)$$

$$= \frac{1}{M} \sum_{j=1}^M A(\mu_{*j}, \sigma_{*j}) - \frac{\sigma_*}{\sqrt{\pi}} - \frac{\frac{1}{M} \sum_{j=1}^M \sigma_j}{\sqrt{\pi}}. \quad (12)$$

- **Moment-matched Gaussian approximation and mixture:**

$$\hat{R}_{\text{Exc}}^{3,2} = \hat{R}_{\text{Exc}}^{2,3} = d(\hat{P}_*, \hat{P}_{\text{ens}}) = \text{CRPS}(\hat{P}_*, \hat{P}_{\text{ens}}) - H(\hat{P}_{\text{ens}}) \quad (13)$$

$$= \frac{1}{M} \sum_{j=1}^M A(\mu_{*j}, \sigma_{*j}) - \frac{\sigma_*}{\sqrt{\pi}} - H(\hat{P}_{\text{ens}}). \quad (14)$$

We notice a peculiar relation:  $2\hat{R}_{\text{Exc}}^{2,1} = \hat{R}_{\text{Exc}}^{1,1}$ . Similar results were obtained for a score with symmetric divergence in (Kotelevskii et al., 2025) in the case of classification.

### 864 A.2.2 LOGARITHMIC SCORE

865 The logarithmic score is a proper scoring rule for both discrete and continuous distributions, which is  
866 defined as follows:

$$867 \text{LS}(\hat{P}, y) = -\log \hat{p}(y), \quad (15)$$

868 where  $\hat{p}(y)$  is the probability density function of the distribution  $\hat{P}$  evaluated at  $y$ . This is an example  
869 of a *local scoring rule*: it depends only on the density value at the point  $y$  and does not take into  
870 account the entire distribution.

871 The entropy function for the logarithmic score is the negative Shannon entropy:

$$872 H(P) = \text{LS}(P, P) = \int_{\mathbb{R}} \text{LS}(P, y) p(y) dy = - \int_{\mathbb{R}} p(y) \log p(y) dy. \quad (16)$$

873 The divergence function of the logarithmic score is the well-known Kullback-Leibler divergence:

$$874 d(\hat{P}, P) = \int_{\mathbb{R}} [-\log \hat{p}(y)] p(y) dy - \int_{\mathbb{R}} [-\log p(y)] p(y) dy = \int_{\mathbb{R}} p(y) \log \frac{p(y)}{\hat{p}(y)} dy = D_{\text{KL}}(P \parallel \hat{P}). \quad (17)$$

875 KL-divergence between two Gaussian mixtures can not be expressed in a closed form, but it can be  
876 approximated with MC methods or bounded, see (Hershey & Olsen, 2007; Durrieu et al., 2012). In  
877 particular, the work (Durrieu et al., 2012) provides lower and upper bounds for the KL divergence  
878 between two Gaussian mixtures using variational approximations.

879 Detailed derivations of the quantities presented in the following are deferred to Section B.2.

880 **Bayes risk.** For the Gaussian ensemble model, we can derive the following approximations of the  
881 Bayes risk for the log score:

- 882 •  $\hat{R}_{\text{Bayes}}^1 = \frac{1}{M} \sum_{i=1}^M H(\hat{P}_i) = \frac{1}{2M} \sum_{i=1}^M \log(2\pi e \sigma_i^2)$ .
- 883 •  $\hat{R}_{\text{Bayes}}^2 = H(\hat{P}_{\text{ens}}) = - \int_{\mathbb{R}} \log \left( \frac{1}{M} \sum_{i=1}^M \hat{p}_i(y) \right) \cdot \left( \frac{1}{M} \sum_{i=1}^M \hat{p}_i(y) \right) dy$  – this expression is  
884 known to have no closed form. Using Jensen’s inequality, we can derive a lower bound:

$$885 \hat{R}_{\text{Bayes}}^2 \geq \frac{1}{M} \sum_{i=1}^M H(\hat{P}_i) = \hat{R}_{\text{Bayes}}^1.$$

- 886 •  $\hat{R}_{\text{Bayes}}^3 = H(\hat{P}_*) = \frac{1}{2} \log(2\pi e \sigma_*^2)$ .

887 **Excess risk.** In the case of the Gaussian ensemble, we can use the following approximations:

- 888 • **Bayesian averaging for both:**

$$889 \hat{R}_{\text{Exc}}^{1,1} = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M d(\hat{P}_i, \hat{P}_j) = \frac{1}{2M^2} \sum_{i=1}^M \sum_{j=1}^M \left[ \frac{\sigma_j^2 + (\mu_i - \mu_j)^2}{\sigma_i^2} - 1 \right]. \quad (18)$$

- 890 • **Gaussian mixture and Bayesian averaging:**

$$891 \hat{R}_{\text{Exc}}^{2,1} = \frac{1}{M} \sum_{i=1}^M d(\hat{P}_{\text{ens}}, \hat{P}_i) = \frac{1}{M} \sum_{i=1}^M D_{\text{KL}}(\hat{P}_i \parallel \hat{P}_{\text{ens}}), \quad (19)$$

$$892 \hat{R}_{\text{Exc}}^{1,2} = \frac{1}{M} \sum_{i=1}^M d(\hat{P}_i, \hat{P}_{\text{ens}}) = \frac{1}{M} \sum_{i=1}^M D_{\text{KL}}(\hat{P}_{\text{ens}} \parallel \hat{P}_i). \quad (20)$$

893 There is no nice analytical form for those since  $D_{\text{KL}}$  contains the log of the mixture density.  
894 The expression for  $\hat{R}_{\text{Exc}}^{2,1}$  was used for discrete distributions in (Lakshminarayanan et al.,  
895 2017) as “disagreement”.

- **Moment-matched Gaussian approximation and Bayesian averaging:**

$$\hat{R}_{\text{Exc}}^{3,1} = \frac{1}{M} \sum_{i=1}^M d(\hat{P}_*, \hat{P}_i) = \frac{1}{2} \left[ \log(\sigma_*^2) - \frac{1}{M} \sum_{i=1}^M \log(\sigma_i^2) \right]. \quad (21)$$

Another, more simple approximation of the mixture is  $\bar{P}_* = \mathcal{N}(\mu_*, \bar{\sigma}_*^2)$ , where  $\bar{\sigma}_*^2 = \frac{1}{M} \sum_{i=1}^M \sigma_i^2$ . With this, we recover the approximation from (Bülte et al., 2025):

$$\bar{R}_{\text{Exc}}^{3,1} = \frac{1}{2} \left[ \log(\bar{\sigma}_*^2) - \frac{1}{M} \sum_{i=1}^M \log(\sigma_i^2) + \frac{1}{\bar{\sigma}_*^2 M} \sum_{i=1}^M (\mu_i - \mu_*)^2 \right].$$

- **Moment-matched Gaussian approximation and mixture:**

$$\hat{R}_{\text{Exc}}^{3,2} = D_{\text{KL}}(\hat{P}_{\text{ens}} \| \hat{P}_*). \quad (22)$$

### A.2.3 QUADRATIC SCORE

In our results below, we will use the following notation:

$$\mathcal{N}(a | b, s^2) = \frac{1}{\sqrt{2\pi}s^2} \exp\left(-\frac{(a-b)^2}{2s^2}\right) = \frac{1}{\sqrt{\sigma_i^2 + \sigma_j^2}} \phi\left(\frac{\mu_i - \mu_j}{\sqrt{\sigma_i^2 + \sigma_j^2}}\right),$$

see Appendix B.3 for more details and derivations of the below quantities.

**Bayes risk.** For the Gaussian ensemble model, we can derive the following approximations of the Bayes risk for the quadratic score:

- $\hat{R}_{\text{Bayes}}^1 = \frac{1}{M} \sum_{i=1}^M H(\hat{P}_i) = -\frac{1}{2M\sqrt{\pi}} \sum_{i=1}^M \frac{1}{\sigma_i}$ .
- $\hat{R}_{\text{Bayes}}^2 = H(\hat{P}_{\text{ens}}) = -\int_{\mathbb{R}} \left(\frac{1}{M} \sum_{i=1}^M \hat{p}_i(y)\right)^2 dy = \frac{1}{2M} \sum_{i=1}^M H(\hat{P}_i) + \frac{1}{2M^2} \sum_{i=1}^M \sum_{j=1}^M \text{QS}(\hat{P}_i, \hat{P}_j)$ .  
Alternatively, we can express the entropy of the quadratic score for a Gaussian mixture as follows:  $\hat{R}_{\text{Bayes}}^2 = -\frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \mathcal{N}(\mu_i | \mu_j, \sigma_i^2 + \sigma_j^2)$ .
- $\hat{R}_{\text{Bayes}}^3 = H(\hat{P}_*) = -\frac{1}{2\sqrt{\pi}\sigma_*^2}$ .

**Excess risk.** The divergence is symmetric for the quadratic score, so the situation is analogous to the case of CRPS. We have two different approximations of Excess risk:

- **Expected Pairwise Bregman Divergence (EPBD):** if we employ Bayesian approach for estimation both  $P$  and  $\hat{P}$ , we get the following approximation:  $\hat{R}_{\text{Exc}}^{1,1} = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M d(\hat{P}_i, \hat{P}_j)$ , where closed form of  $d(\hat{P}_i, \hat{P}_j)$  is given by (see equation (71)):

$$d(\hat{P}_i, \hat{P}_j) = \frac{1}{2\sqrt{\pi}\sigma_i} + \frac{1}{2\sqrt{\pi}\sigma_j} - 2\mathcal{N}(\mu_i | \mu_j, \sigma_i^2 + \sigma_j^2).$$

We combine these results to get the final expression:

$$\hat{R}_{\text{Exc}}^{1,1} = \underbrace{\frac{1}{M\sqrt{\pi}} \sum_{i=1}^M \frac{1}{\sigma_i}}_{-2\hat{R}_{\text{Bayes}}^1} - \frac{2}{M^2} \sum_{i=1}^M \sum_{j=1}^M \mathcal{N}(\mu_i | \mu_j, \sigma_i^2 + \sigma_j^2).$$

Notice that the first part equal  $-2\hat{R}_{\text{Bayes}}^1$ .

- 972 • **Bregman Information (BI):** here we combine Bayesian averaging and the mixture distri-  
973 bution. Using equation (75) for the divergence between the mixture and its component we  
974 get:

$$\begin{aligned}
 \widehat{R}_{\text{Exc}}^{2,1} &= \widehat{R}_{\text{Exc}}^{1,2} = \frac{1}{M} \sum_{i=1}^M d(\widehat{P}_i, \widehat{P}_{\text{ens}}) \\
 &= \frac{1}{M} \sum_{i=1}^M \left[ -\frac{2}{M} \left( \sum_{j=1}^M \mathcal{N}(\mu_i | \mu_j, \sigma_i^2 + \sigma_j^2) \right) - H(\widehat{P}_i) - H(\widehat{P}_{\text{ens}}) \right] \\
 &= \underbrace{\frac{1}{2M\sqrt{\pi}} \sum_{i=1}^M \frac{1}{\sigma_i}}_{-\widehat{R}_{\text{Bayes}}^1} - \frac{2}{M^2} \sum_{i=1}^M \sum_{j=1}^M \mathcal{N}(\mu_i | \mu_j, \sigma_i^2 + \sigma_j^2) - \widehat{R}_{\text{Bayes}}^2 = \frac{1}{2} \widehat{R}_{\text{Exc}}^{1,1}.
 \end{aligned}$$

- 986 • **Moment-matched Gaussian approximation and Bayesian averaging:**

$$\begin{aligned}
 \widehat{R}_{\text{Exc}}^{3,1} &= \widehat{R}_{\text{Exc}}^{1,3} = \frac{1}{M} \sum_{i=1}^M d(\widehat{P}_*, \widehat{P}_i) \\
 &= \frac{1}{M} \sum_{i=1}^M \left[ -H(\widehat{P}_i) - H(\widehat{P}_*) - 2\mathcal{N}(\mu_* | \mu_i, \sigma_*^2 + \sigma_i^2) \right] \\
 &= \underbrace{\frac{1}{2M\sqrt{\pi}} \sum_{i=1}^M \frac{1}{\sigma_i}}_{-\widehat{R}_{\text{Bayes}}^1} - \frac{2}{M} \sum_{i=1}^M \mathcal{N}(\mu_* | \mu_i, \sigma_*^2 + \sigma_i^2) + \underbrace{\frac{1}{2\sqrt{\pi}\sigma_*^2}}_{-\widehat{R}_{\text{Bayes}}^3}.
 \end{aligned}$$

- 999 • **Moment-matched Gaussian approximation and mixture:**

$$\widehat{R}_{\text{Exc}}^{3,2} = \widehat{R}_{\text{Exc}}^{2,3} = d(\widehat{P}_*, \widehat{P}_{\text{ens}}) = -\frac{2}{M} \left( \sum_{j=1}^M \mathcal{N}(\mu_* | \mu_j, \sigma_*^2 + \sigma_j^2) \right) - \underbrace{H(\widehat{P}_*)}_{\widehat{R}_{\text{Bayes}}^3} - \underbrace{H(\widehat{P}_{\text{ens}})}_{\widehat{R}_{\text{Bayes}}^2}.$$

#### 1004 A.2.4 SE SCORE

1006 We can use the same three approaches to estimate the risk components for the SE score as we did  
1007 for other proper scoring rules. Due to the symmetry of the divergence, the central label and central  
1008 prediction coincide again. They are equal to the full mixture distribution, which we again denote  
1009 as  $\widehat{P}_{\text{ens}}$ . The mean and variance of this distribution are denoted as  $\mu_*$  and  $\sigma_*^2$ . We have previously  
1010 introduced them to construct our Gaussian approximation in equation (5). Here they are again for  
1011 clarity:

$$\mathbb{E}[\widehat{P}_{\text{ens}}] = \mu_* = \frac{1}{M} \sum_{i=1}^M \mu_i, \quad \text{Var}[\widehat{P}_{\text{ens}}] = \sigma_*^2 = \frac{1}{M} \sum_{i=1}^M (\sigma_i^2 + \mu_i^2) - \mu_*^2.$$

1015 For a more detailed treatment of the SE score, please refer to Section C.

#### 1016 Bayes risk.

- 1018 • **Bayesian averaging:**  $\widehat{R}_{\text{Bayes}}^1 = \frac{1}{M} \sum_{i=1}^M H_{\text{SE}}(\widehat{P}_i) = \frac{1}{M} \sum_{i=1}^M \text{Var}_{Y \sim \widehat{P}_i}[Y] =$   
1019  $\frac{1}{M} \sum_{i=1}^M \sigma_i^2.$
- 1021 • **Posterior predictive:**  $\widehat{R}_{\text{Bayes}}^2 = H_{\text{SE}}(\widehat{P}_{\text{ens}}) = \text{Var}_{Y \sim \widehat{P}_{\text{ens}}}[Y] = \sigma_*^2 = \frac{1}{M} \sum_{i=1}^M (\sigma_i^2 + \mu_i^2) -$   
1022  $\mu_*^2.$
- 1023 • **Gaussian approximation:**  $\widehat{R}_{\text{Bayes}}^3 = H_{\text{SE}}(\widehat{P}_*) = \text{Var}_{Y \sim \widehat{P}_*}[Y] = \sigma_*^2 = \widehat{R}_{\text{Bayes}}^2.$

1025 **Excess risk.** For the symmetric divergence of SE we get:

- 1026 • **Bayesian averaging:**  $\widehat{R}_{\text{Exc}}^{1,1} = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M d_{\text{SE}}(\widehat{P}_i, \widehat{P}_j) = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M (\mu_j -$   
 1027  $\mu_i)^2 = 2\widehat{\text{Var}}[\mu_i]$ .  
 1028
- 1029 • **Posterior predictive and Bayesian averaging:**  $\widehat{R}_{\text{Exc}}^{2,1} = \widehat{R}_{\text{Exc}}^{1,2} = \frac{1}{M} \sum_{i=1}^M d_{\text{SE}}(\widehat{P}_i, \widehat{P}_{\text{ens}}) =$   
 1030  $\frac{1}{M} \sum_{i=1}^M (\mathbb{E}_{X \sim \widehat{P}_i}[X] - \mathbb{E}_{Y \sim \widehat{P}_{\text{ens}}}[Y])^2 = \frac{1}{M} \sum_{i=1}^M (\mu_i - \mu_*)^2 = \widehat{\text{Var}}[\mu_i]$ .  
 1031
- 1032 • **Moment-matched Gaussian approximation and Bayesian averaging:** since SE score  
 1033 only looks at the mean of the distribution, using this Gaussian approximation will give the  
 1034 same result as the previous section.

$$\widehat{R}_{\text{Exc}}^{3,1} = \widehat{R}_{\text{Exc}}^{1,3} = \frac{1}{M} \sum_{i=1}^M (\mu_i - \mu_*)^2 = \widehat{\text{Var}}[\mu_i].$$

- 1038 • **Moment-matched Gaussian approximation and posterior predictive:** the means of these  
 1039 distributions coincide so we get a zero.

$$\widehat{R}_{\text{Exc}}^{3,2} = \widehat{R}_{\text{Exc}}^{2,3} = 0.$$

## 1043 B SCORE COMPUTATION AND TOOLS

### 1045 B.1 CRPS

1047 Continuous ranked probability score has multiple equivalent representations:

$$1048 \text{CRPS}(P, y) = \int_{\mathbb{R}} (F_P(t) - \mathbb{I}\{y \leq t\})^2 dt. \quad (23)$$

$$1051 \text{CRPS}(P, y) = \mathbb{E}_{X \sim P} |X - y| - \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'|. \quad (24)$$

$$1053 \text{CRPS}(P, y) = \int |x - y| dP(x) - \frac{1}{2} \int \int |x - x'| dP(x) dP(x'). \quad (25)$$

1055 Expected CRPS between two distributions  $P$  and  $Q$  can be expressed as:

$$1056 \text{CRPS}(P, Q) = \mathbb{E}_{X \sim P, Y \sim Q} |X - Y| - \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'|. \quad (26)$$

#### 1060 B.1.1 ENTROPY FUNCTION OF CRPS

$$1062 H(P) = \int \text{CRPS}(P, y) dP(y) = \int \left[ \mathbb{E}_{X \sim P} |X - y| - \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'| \right] dP(y)$$

$$1064 = \int \mathbb{E}_{X \sim P} |X - y| dP(y) - \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'| = \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'|.$$

#### 1067 B.1.2 DIVERGENCE FUNCTION OF CRPS

$$1068 d(P, Q) = \int \text{CRPS}(P, y) dQ(y) - H(Q)$$

$$1070 = \int \left[ \mathbb{E}_{X \sim P} |X - y| - \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'| \right] dQ(y) - H(Q)$$

$$1072 = \mathbb{E}_{X \sim P, Y \sim Q} |X - Y| - \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'| - \frac{1}{2} \mathbb{E}_{Y, Y' \sim Q} |Y - Y'|$$

$$1074 = \mathbb{E}_{X \sim P, Y \sim Q} |X - Y| - H(P) - H(Q).$$

1077 Another expression for the divergence function of CRPS is given by:

$$1078 d(P, Q) = \int_{\mathbb{R}} (F_P(y) - F_Q(y))^2 dy. \quad (27)$$

### B.1.3 CRPS FOR GAUSSIANS

**Notation.** Some commonly used notation:

$$\phi(t) = \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} \text{ is the standard normal density function.} \quad (28)$$

$$\Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} \text{ is the standard normal cumulative distribution function.} \quad (29)$$

We follow (Villegz, 2017) and also introduce the following notation:

$$\mu_{ij} := \mu_i - \mu_j, \quad (30)$$

$$\sigma_{ij} := \sqrt{\sigma_i^2 + \sigma_j^2}, \quad (31)$$

$$A(\mu, \sigma) := 2\sigma\phi\left(\frac{\mu}{\sigma}\right) + \mu \left[2\Phi\left(\frac{\mu}{\sigma}\right) - 1\right]. \quad (32)$$

We note that for  $X \sim \mathcal{N}(\mu, \sigma)$  it can be shown that  $\mathbb{E}[|X|] = A(\mu, \sigma)$ .

**CRPS for a Gaussian.**

$$\text{CRPS}(\mathcal{N}(\mu, \sigma), y) = \sigma \left( 2\phi(z) + z(2\Phi(z) - 1) - \frac{1}{\sqrt{\pi}} \right), \quad (33)$$

where  $z = \frac{y-\mu}{\sigma}$ ,  $\phi(z)$  is PDF of the standard normal distribution, and  $\Phi(z)$  is CDF of the standard normal distribution

*Proof.*

$$\begin{aligned} \text{CRPS}(\mathcal{N}(\mu, \sigma), y) &= \int_{-\infty}^{+\infty} (F_{\mathcal{N}}(t) - \mathbb{I}\{y \leq t\})^2 dt \\ &= \int_{-\infty}^y (F_{\mathcal{N}}(t) - \mathbb{I}\{y \leq t\})^2 dt + \int_y^{+\infty} (F_{\mathcal{N}}(t) - \mathbb{I}\{y \leq t\})^2 dt \\ &= \int_{-\infty}^y (F_{\mathcal{N}}(t) - 0)^2 dt + \int_y^{+\infty} (F_{\mathcal{N}}(t) - 1)^2 dt \\ &= \int_{-\infty}^y F_{\mathcal{N}}^2(t) dt + \int_y^{+\infty} (1 - F_{\mathcal{N}}(t))^2 dt = I_1 + I_2. \end{aligned}$$

First we compute

$$I_1 = \int_{-\infty}^y F_{\mathcal{N}}^2(t) dt = \int_{-\infty}^y \Phi^2\left(\frac{t-\mu}{\sigma}\right) dt = \sigma \int_{-\infty}^{\frac{y-\mu}{\sigma}} \Phi^2(z) dz = \sigma \int_{-\infty}^{z_y} \Phi^2(z) dz,$$

where we used the substitution  $z = \frac{t-\mu}{\sigma}$  and define  $z_y = \frac{y-\mu}{\sigma}$ . For the second part, we get

$$\begin{aligned} I_2 &= \int_y^{+\infty} (1 - F_{\mathcal{N}}(t))^2 dt = \int_y^{+\infty} \left(1 - \Phi\left(\frac{t-\mu}{\sigma}\right)\right)^2 dt \\ &= \sigma \int_{\frac{y-\mu}{\sigma}}^{+\infty} (1 - \Phi(z))^2 dz = \sigma \int_{z_y}^{+\infty} (1 - \Phi(z))^2 dz. \end{aligned}$$

Let us calculate the first part:

$$\begin{aligned} I_1 &= \sigma \int_{-\infty}^{z_y} \Phi^2(z) dz = \left[ \begin{array}{l} u = \Phi(z)^2 \\ du = 2\Phi(z)\phi(z) dz \end{array} \quad \begin{array}{l} dv = dz \\ v = z \end{array} \right] = \sigma \left[ z\Phi^2(z) \Big|_{-\infty}^{z_y} - \int_{-\infty}^{z_y} 2z\Phi(z)\phi(z) dz \right] \\ &= \sigma \left[ z_y\Phi^2(z_y) - 2 \int_{-\infty}^{z_y} z\Phi(z)\phi(z) dz \right] = \sigma \left[ z_y\Phi^2(z_y) + 2 \int_{-\infty}^{z_y} z\Phi(z) d\phi(z) \right], \end{aligned}$$

where we have used the fact that  $(\phi(z))'_z = -z\phi(z)$ .

Next,

$$\begin{aligned}
I_2 &= \sigma \int_{z_y}^{+\infty} (1 - \Phi(z))^2 dz = \left[ \begin{array}{ll} u = (1 - \Phi(z))^2 & dv = dz \\ du = 2(\Phi(z) - 1)\phi(z)dz & v = z \end{array} \right] \\
&= \sigma \left[ z(1 - \Phi(z))^2 \Big|_{z_y}^{+\infty} - 2 \int_{z_y}^{+\infty} z(\Phi(z) - 1)\phi(z) dz \right] \\
&= \sigma \left[ -z_y(1 - \Phi(z_y))^2 + 2 \int_{z_y}^{+\infty} (\Phi(z) - 1)d\phi(z) \right]. \tag{34}
\end{aligned}$$

Finally,

$$\begin{aligned}
\text{CRPS}(\mathcal{N}(\mu, \sigma), y) &= I_1 + I_2 = \sigma \left[ z_y(2\Phi(z_y) - 1) + 2 \int_{-\infty}^{+\infty} \Phi(z) d\phi(z) - 2 \int_{z_y}^{+\infty} d\phi(z) \right] \\
&= \sigma [z_y(2\Phi(z_y) - 1) + 2I_3 + 2\phi(z_y)] \\
&= \sigma \left[ z_y(2\Phi(z_y) - 1) - \frac{1}{\sqrt{\pi}} + 2\phi(z_y) \right]. \tag{35}
\end{aligned}$$

We used the following identity:

$$I_3 = \int_{-\infty}^{+\infty} \Phi(z) d\phi(z) = \Phi(z)\phi(z) \Big|_{-\infty}^{+\infty} - \int_{-\infty}^{+\infty} \phi^2(z) dz = -\frac{1}{2\sqrt{\pi}}.$$

□

**Entropy of CRPS for a Gaussian.** If we use the alternative representation of CRPS, we can express it as the expected absolute value of a centered Gaussian:

$$\begin{aligned}
X, X' \sim \mathcal{N}(\mu, \sigma^2) &\implies X - X' \sim \mathcal{N}(0, 2\sigma^2) \\
&\implies \mathbb{E}_{X, X' \sim P} |X - X'| = A(0, \sqrt{2}\sigma) = \sqrt{2}\sigma \sqrt{\frac{2}{\pi}} = \frac{2\sigma}{\sqrt{\pi}}. \tag{36}
\end{aligned}$$

$$H(\mathcal{N}(\mu, \sigma^2)) = \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'| = \frac{1}{2} \frac{2\sigma}{\sqrt{\pi}} = \frac{\sigma}{\sqrt{\pi}}. \tag{37}$$

**Divergence function of CRPS for two Gaussians.** Let  $P = \mathcal{N}(\mu_i, \sigma_i^2)$  and  $Q = \mathcal{N}(\mu_j, \sigma_j^2)$  be two Gaussian distributions. Recall that the divergence function of CRPS is given by:

$$d(P, Q) = \mathbb{E}_{X \sim P, Y \sim Q} |X - Y| - \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'| - \frac{1}{2} \mathbb{E}_{Y, Y' \sim Q} |Y - Y'|. \tag{38}$$

Since  $X$  and  $Y$  are independent,  $X - Y \sim \mathcal{N}(\mu_i - \mu_j, \sigma_i^2 + \sigma_j^2)$ . Using our notation, we get:

$$\mathbb{E}_{X \sim P, Y \sim Q} |X - Y| = A(\mu_{ij}, \sigma_{ij}). \tag{39}$$

We also note that  $H(P) = \frac{\sigma_i}{\sqrt{\pi}}$  and  $H(Q) = \frac{\sigma_j}{\sqrt{\pi}}$ . Combining all these results, we can express the divergence function as follows:

$$d(P, Q) = A(\mu_{ij}, \sigma_{ij}) - \frac{\sigma_i + \sigma_j}{\sqrt{\pi}}. \tag{40}$$

**Expected CRPS between two Gaussians.**

$$\text{CRPS}(P, Q) = \mathbb{E}_{X \sim P, Y \sim Q} |X - Y| - \frac{1}{2} \mathbb{E}_{X, X' \sim P} |X - X'|. \tag{41}$$

Let  $P = \mathcal{N}(\mu_i, \sigma_i^2)$  and  $Q = \mathcal{N}(\mu_j, \sigma_j^2)$  be two Gaussian distributions. Then the expected CRPS between these two distributions can be expressed as:

$$\text{CRPS}(\mathcal{N}(\mu_i, \sigma_i), \mathcal{N}(\mu_j, \sigma_j)) = A(\mu_{ij}, \sigma_{ij}) - \frac{\sigma_i}{\sqrt{\pi}}. \tag{42}$$

**Entropy of CRPS for Gaussian mixture.**

$$H(\hat{P}_{\text{ens}}) = \frac{1}{2} \mathbb{E}_{X, X' \sim \hat{P}_{\text{ens}}} |X - X'| = \frac{1}{2} \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \mathbb{E}_{X \sim \hat{P}_i, X' \sim \hat{P}_j} |X - X'| = \frac{1}{2} \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M A(\mu_{ij}, \sigma_{ij}). \quad (43)$$

**Expected CRPS between a Gaussian and a Gaussian mixture.**

$$\text{CRPS}(\hat{P}_i, \hat{P}_{\text{ens}}) = \mathbb{E}_{X \sim \hat{P}_i, Y \sim \hat{P}_{\text{ens}}} |X - Y| - \frac{1}{2} \mathbb{E}_{X, X' \sim \hat{P}_i} |X - X'| = \frac{1}{M} \sum_{j=1}^M A(\mu_{ij}, \sigma_{ij}) - \frac{\sigma_i}{\sqrt{\pi}}. \quad (44)$$

**Divergence between a Gaussian and a Gaussian mixture.**

$$d(\hat{P}_i, \hat{P}_{\text{ens}}) = \text{CRPS}(\hat{P}_i, \hat{P}_{\text{ens}}) - H(\hat{P}_{\text{ens}}) = \frac{1}{M} \sum_{j=1}^M A(\mu_{ij}, \sigma_{ij}) - \frac{\sigma_i}{\sqrt{\pi}} - \frac{1}{2} \frac{1}{M^2} \sum_{l=1}^M \sum_{j=1}^M A(\mu_{lj}, \sigma_{lj}). \quad (45)$$

**B.2 LOG SCORE****B.2.1 ENTROPY OF LOG SCORE**

Log score is associated with Shannon entropy:

$$H(P) = \mathbb{E}_{Y \sim P} [-\log p(Y)] = - \int_{\mathbb{R}} p(y) \log p(y) dy. \quad (46)$$

**B.2.2 DIVERGENCE OF LOG SCORE**

Divergence function is the Kullback-Leibler divergence:

$$d(\hat{P}, P) = D_{\text{KL}}(P \parallel \hat{P}) = \int_{\mathbb{R}} p(y) \log \frac{p(y)}{\hat{p}(y)} dy. \quad (47)$$

**B.2.3 LOG SCORE FOR GAUSSIANS**

**Log score entropy for a Gaussian.** Logarithm of the Gaussian density is a quadratic function, so we can easily express the entropy using the expression for Gaussian:

$$Y \sim P = \mathcal{N}(\mu, \sigma) \implies p(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \implies \text{LS}(P, y) = -\log p(y) = \frac{1}{2} \log(2\pi\sigma^2) + \frac{(y-\mu)^2}{2\sigma^2}. \quad (48)$$

$$\begin{aligned} H(\mathcal{N}(\mu, \sigma)) &= \mathbb{E}_{Y \sim \mathcal{N}(\mu, \sigma)} [-\log p(Y)] = \mathbb{E}_{Y \sim \mathcal{N}(\mu, \sigma)} \left[ \frac{1}{2} \log(2\pi\sigma^2) + \frac{(y-\mu)^2}{2\sigma^2} \right] \\ &= \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \mathbb{E}_{Y \sim \mathcal{N}(\mu, \sigma)} [(y-\mu)^2] = \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} = \frac{1}{2} \log(2\pi e\sigma^2). \end{aligned}$$

**Log score divergence for two Gaussians.** Here, the derivation is very similar, using the variance/second moment formulas.

$$P = \mathcal{N}(\mu_p, \sigma_p^2), Q = \mathcal{N}(\mu_q, \sigma_q^2) \implies p(y) = \frac{1}{\sqrt{2\pi}\sigma_p} e^{-\frac{(y-\mu_p)^2}{2\sigma_p^2}}, q(y) = \frac{1}{\sqrt{2\pi}\sigma_q} e^{-\frac{(y-\mu_q)^2}{2\sigma_q^2}}. \quad (49)$$

$$\log p(y) = -\frac{1}{2} \log(2\pi\sigma_p^2) - \frac{(y-\mu_p)^2}{2\sigma_p^2}, \log q(y) = -\frac{1}{2} \log(2\pi\sigma_q^2) - \frac{(y-\mu_q)^2}{2\sigma_q^2}. \quad (50)$$

$$\begin{aligned}
d(P, Q) &= \mathbb{E}_{Y \sim Q} [-\log p(Y)] - H(Q) = D_{\text{KL}}(Q \| P) \\
&= \mathbb{E}_{Y \sim Q} [\log q(Y) - \log p(Y)] = \mathbb{E}_{Y \sim Q} \left[ \frac{1}{2} \log(2\pi\sigma_p^2) + \frac{(y - \mu_p)^2}{2\sigma_p^2} - \frac{1}{2} \log(2\pi\sigma_q^2) - \frac{(y - \mu_q)^2}{2\sigma_q^2} \right] \\
&= \frac{1}{2} \left[ \log \frac{\sigma_p^2}{\sigma_q^2} - 1 + \frac{\sigma_q^2 + (\mu_p - \mu_q)^2}{\sigma_p^2} \right].
\end{aligned}$$

### Expected Log score between two Gaussians.

$$\begin{aligned}
\text{LS}(P, Q) &= d(P, Q) + H(Q) = \frac{1}{2} \left[ \log \frac{\sigma_p^2}{\sigma_q^2} - 1 + \frac{\sigma_q^2 + (\mu_p - \mu_q)^2}{\sigma_p^2} \right] + \frac{1}{2} \log(2\pi e \sigma_q^2) \\
&= \frac{1}{2} \left[ \log(2\pi\sigma_p^2) + \frac{\sigma_q^2 + (\mu_p - \mu_q)^2}{\sigma_p^2} \right]. \tag{51}
\end{aligned}$$

### Pairwise divergences of a Gaussian ensemble.

$$\begin{aligned}
\hat{R}_{\text{Exc}}^{1,1} &= \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M d(\hat{P}_i, \hat{P}_j) = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M D_{\text{KL}}(\hat{P}_j, \hat{P}_i) \\
&= \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \frac{1}{2} \left[ \log \frac{\sigma_i^2}{\sigma_j^2} - 1 + \frac{\sigma_j^2 + (\mu_i - \mu_j)^2}{\sigma_i^2} \right] \\
&= \frac{1}{2M^2} \sum_{i=1}^M \sum_{j=1}^M \left[ \frac{\sigma_j^2 + (\mu_i - \mu_j)^2}{\sigma_i^2} - 1 \right], \tag{52}
\end{aligned}$$

because sum of all pairwise differences of logarithms is 0.

**Divergence between a Gaussian mixture and one of its components.** For the log score the divergence function is KL-divergence, it is not symmetric. This provides two possible ways to approximate the excess risk: assume the true distribution to be the full mixture and average over the components or the other way around. In the notation of (Kotelevskii et al., 2025) they correspond to Bregman Information (BI) and Reverse Bregman Information (RBI).

$$\begin{aligned}
d(\hat{P}_{\text{ens}}, \hat{P}_i) &= D_{\text{KL}}(\hat{P}_i \| \hat{P}_{\text{ens}}) \\
&= \text{LS}(\hat{P}_{\text{ens}}, \hat{P}_i) - H(\hat{P}_i) = - \int_{\mathbb{R}} \hat{p}_i(y) \log(\hat{p}_{\text{ens}}(y)) dy - \frac{1}{2} \log(2\pi e \sigma_i^2). \\
d(\hat{P}_i, \hat{P}_{\text{ens}}) &= D_{\text{KL}}(\hat{P}_{\text{ens}} \| \hat{P}_i) = \text{LS}(\hat{P}_i, \hat{P}_{\text{ens}}) - H(\hat{P}_{\text{ens}}) \\
&= - \int_{\mathbb{R}} \hat{p}_{\text{ens}}(y) \log \hat{p}_i(y) dy - H(\hat{P}_{\text{ens}}) = \frac{1}{M} \sum_{j=1}^M \text{LS}(\hat{P}_i, \hat{P}_j) - H(\hat{P}_{\text{ens}}) \\
&= \frac{1}{M} \sum_{j=1}^M \text{LS}(\hat{P}_i, \hat{P}_j) - H(\hat{P}_{\text{ens}}).
\end{aligned}$$

There is no more explicit analytical formula due to appearance of logarithm of the mixture density, which is a sum of the component densities. In the discrete case there are expressions involving the *LogSumExp* function.

### Divergence between a Gaussian ensemble and its Gaussian approximation.

$$\hat{P}_* = \mathcal{N}(\mu_*, \sigma_*), \mu_* = \frac{1}{M} \sum_{i=1}^M \mu_i, \sigma_*^2 = \frac{1}{M} \sum_{i=1}^M (\sigma_i^2 + \mu_i^2) - \mu_*^2.$$

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307

$$\begin{aligned}\widehat{R}_{\text{Exc}}^{3a,1} &= \frac{1}{M} \sum_{i=1}^M d(\widehat{P}_*, \widehat{P}_i) = \frac{1}{M} \sum_{i=1}^M \frac{1}{2} \left[ \log \frac{\sigma_*^2}{\sigma_i^2} - 1 + \frac{\sigma_i^2 + (\mu_* - \mu_i)^2}{\sigma_*^2} \right] \\ &= \frac{1}{2} \left[ \log(\sigma_*^2) - \frac{1}{M} \sum_{i=1}^M \log(\sigma_i^2) - 1 + \frac{1}{\sigma_*^2} \frac{1}{M} \sum_{i=1}^M \sigma_i^2 + \frac{(\mu_* - \mu_i)^2}{\sigma_*^2} \right] \\ &= \frac{1}{2} \left[ \log(\sigma_*^2) - \frac{1}{M} \sum_{i=1}^M \log(\sigma_i^2) \right].\end{aligned}$$

1308  
1309  
1310

**Divergence between a Gaussian ensemble and mean Gaussian approximation.**  $\bar{P}_* = \mathcal{N}(\mu_*, \bar{\sigma}_*^2)$ , where  $\mu_* = \frac{1}{M} \sum_{i=1}^M \mu_i$  and  $\bar{\sigma}_*^2 = \frac{1}{M} \sum_{i=1}^M \sigma_i^2$ .

1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322

$$\begin{aligned}\widehat{R}_{\text{Exc}}^{3b,1} &= \frac{1}{M} \sum_{i=1}^M d(\widehat{P}, \widehat{P}_i) = \frac{1}{M} \sum_{i=1}^M \frac{1}{2} \left[ \log \frac{\bar{\sigma}_*^2}{\sigma_i^2} - 1 + \frac{\sigma_i^2 + (\mu_* - \mu_i)^2}{\bar{\sigma}_*^2} \right] \\ &= \frac{1}{2} \left[ \log(\bar{\sigma}_*^2) - \frac{1}{M} \sum_{i=1}^M \log(\sigma_i^2) - 1 + \frac{1}{\bar{\sigma}_*^2 M} \sum_{i=1}^M \sigma_i^2 + \frac{1}{\bar{\sigma}_*^2 M} \sum_{i=1}^M (\mu_* - \mu_i)^2 \right] \\ &= \frac{1}{2} \left[ \log(\bar{\sigma}_*^2) - \frac{1}{M} \sum_{i=1}^M \log(\sigma_i^2) + \frac{1}{\bar{\sigma}_*^2 M} \sum_{i=1}^M (\mu_* - \mu_i)^2 \right].\end{aligned}$$

1323

### B.3 QUADRATIC SCORE

1324

Quadratic score is a generalization of the Brier score for continuous outcomes:

1325  
1326  
1327

$$\text{QS}(\widehat{P}, y) = -2\widehat{p}(y) + \int_{\mathbb{R}} \widehat{p}(t)^2 dt. \quad (53)$$

1328  
1329  
1330

It can also be viewed as a way to make the value of the probability density at the outcome  $y$  into a proper scoring rule. Simply taking  $S(\widehat{P}, y) = -\widehat{p}(y)$  is not a proper scoring rule.

1331

1332

#### B.3.1 QUADRATIC SCORE ENTROPY

1333

1334

1335

1336

1337

1338

$$\begin{aligned}H(P) &= \mathbb{E}_{Y \sim P} [\text{QS}(P, Y)] = -2\mathbb{E}_{Y \sim P} [p(Y)] + \int_{\mathbb{R}} p(t)^2 dt \\ &= -2 \int_{\mathbb{R}} p(y)p(y) dy + \int_{\mathbb{R}} p(t)^2 dt = - \int_{\mathbb{R}} p(y)^2 dy.\end{aligned} \quad (54)$$

1339

1340

#### B.3.2 QUADRATIC SCORE ENTROPY FOR A GAUSSIAN

1341

1342

1343

1344

1345

$$\begin{aligned}H(\mathcal{N}(\mu, \sigma)) &= - \int_{\mathbb{R}} p(y)^2 dy = - \int_{\mathbb{R}} \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t-\mu)^2}{2\sigma^2}} \right)^2 dt = - \frac{1}{2\pi\sigma^2} \int_{\mathbb{R}} e^{-\frac{(t-\mu)^2}{\sigma^2}} dt \\ &= \left[ u = \frac{t-\mu}{\sigma} \right] = - \frac{1}{2\pi\sigma^2} \sigma \int_{\mathbb{R}} e^{-u^2} du = \frac{\sqrt{\pi}}{2\pi\sigma} = - \frac{1}{2\sqrt{\pi}\sigma}.\end{aligned}$$

1346

1347

1348

1349

#### B.3.3 EXPECTED QUADRATIC SCORE FOR A GAUSSIAN PREDICTION AND GAUSSIAN OUTCOME

$$\mathbb{E}_{Y \sim Q} [\text{QS}(P, Y)] = -2\mathbb{E}_{Y \sim Q} [p(Y)] + \int_{\mathbb{R}} p(t)^2 dt = -2 \int_{\mathbb{R}} p(y)q(y) dy - H(P). \quad (55)$$

We will focus first on computing the integral. Since  $P$  and  $Q$  are Gaussian distributions, we can express the integral as follows:

$$\int_{\mathbb{R}} p(t)q(t)dt = \frac{1}{2\pi\sigma_P\sigma_Q} \int_{\mathbb{R}} e^{-\frac{(t-\mu_P)^2}{2\sigma_P^2} - \frac{(t-\mu_Q)^2}{2\sigma_Q^2}} dt = \frac{1}{2\pi\sigma_P\sigma_Q} \int_{\mathbb{R}} e^{-\frac{1}{2}\left(\frac{(t-\mu_P)^2}{\sigma_P^2} + \frac{(t-\mu_Q)^2}{\sigma_Q^2}\right)} dt. \quad (56)$$

Let us simplify the exponent:

$$\begin{aligned} \frac{(t-\mu_P)^2}{\sigma_P^2} + \frac{(t-\mu_Q)^2}{\sigma_Q^2} &= \frac{t^2}{\sigma_P^2} - 2\frac{t\mu_P}{\sigma_P^2} + \frac{\mu_P^2}{\sigma_P^2} + \frac{t^2}{\sigma_Q^2} - 2\frac{t\mu_Q}{\sigma_Q^2} + \frac{\mu_Q^2}{\sigma_Q^2} \\ &= \left(\frac{1}{\sigma_P^2} + \frac{1}{\sigma_Q^2}\right)t^2 - 2\left(\frac{\mu_P}{\sigma_P^2} + \frac{\mu_Q}{\sigma_Q^2}\right)t + \left(\frac{\mu_P^2}{\sigma_P^2} + \frac{\mu_Q^2}{\sigma_Q^2}\right) = At^2 - 2Bt + C = A\left(t - \frac{B}{A}\right)^2 - \frac{B^2}{A} + C, \end{aligned} \quad (57)$$

where we denote:

$$A = \frac{1}{\sigma_P^2} + \frac{1}{\sigma_Q^2}, \quad B = \frac{\mu_P}{\sigma_P^2} + \frac{\mu_Q}{\sigma_Q^2}, \quad C = \frac{\mu_P^2}{\sigma_P^2} + \frac{\mu_Q^2}{\sigma_Q^2}. \quad (58)$$

Now we have:

$$\int_{\mathbb{R}} e^{-\frac{1}{2}(At^2 - 2Bt + C)} dt = e^{-\frac{1}{2}(C - \frac{B^2}{A})} \int_{\mathbb{R}} e^{-\frac{A}{2}(t - \frac{B}{A})^2} dt. \quad (59)$$

The leftover integral is a Gaussian integral, which can be computed as follows:

$$\int_{\mathbb{R}} e^{-\frac{A}{2}(t - \frac{B}{A})^2} dt = \sqrt{\frac{2\pi}{A}} = \sqrt{2\pi} \frac{\sigma_P\sigma_Q}{\sqrt{\sigma_P^2 + \sigma_Q^2}}. \quad (60)$$

Now, we simplify the exponent of the coefficient in front of that integral:

$$A = \frac{1}{\sigma_P^2} + \frac{1}{\sigma_Q^2} = \frac{\sigma_Q^2 + \sigma_P^2}{\sigma_P^2\sigma_Q^2}, \quad B^2 = \left(\frac{\mu_P}{\sigma_P^2} + \frac{\mu_Q}{\sigma_Q^2}\right)^2 = \frac{(\mu_P\sigma_Q^2 + \mu_Q\sigma_P^2)^2}{\sigma_P^4\sigma_Q^4}, \quad C = \frac{\mu_P^2\sigma_Q^2 + \mu_Q^2\sigma_P^2}{\sigma_P^2\sigma_Q^2}. \quad (61)$$

$$C - \frac{B^2}{A} = \frac{\mu_P^2\sigma_Q^2 + \mu_Q^2\sigma_P^2}{\sigma_P^2\sigma_Q^2} - \frac{(\mu_P\sigma_Q^2 + \mu_Q\sigma_P^2)^2}{\sigma_P^2\sigma_Q^2(\sigma_Q^2 + \sigma_P^2)} = \frac{(\mu_P^2\sigma_Q^2 + \mu_Q^2\sigma_P^2)(\sigma_Q^2 + \sigma_P^2) - (\mu_P\sigma_Q^2 + \mu_Q\sigma_P^2)^2}{\sigma_P^2\sigma_Q^2(\sigma_Q^2 + \sigma_P^2)}. \quad (62)$$

$$\begin{aligned} &(\mu_P^2\sigma_Q^2 + \mu_Q^2\sigma_P^2)(\sigma_Q^2 + \sigma_P^2) - (\mu_P\sigma_Q^2 + \mu_Q\sigma_P^2)^2 \\ &= \cancel{\mu_P^2\sigma_Q^4} + \cancel{\mu_Q^2\sigma_P^4} + \mu_P\sigma_P^2\sigma_Q^2 + \mu_Q\sigma_P^2\sigma_Q^2 - (\cancel{\mu_P^2\sigma_Q^4} + 2\mu_P\mu_Q\sigma_P^2\sigma_Q^2 + \cancel{\mu_Q^2\sigma_P^4}) \\ &= \sigma_P^2\sigma_Q^2(\mu_P - \mu_Q)^2. \end{aligned} \quad (63)$$

Thus, we have:

$$C - \frac{B^2}{A} = \frac{\sigma_P^2\sigma_Q^2(\mu_P - \mu_Q)^2}{\sigma_P^2\sigma_Q^2(\sigma_Q^2 + \sigma_P^2)} = \frac{(\mu_P - \mu_Q)^2}{\sigma_Q^2 + \sigma_P^2}. \quad (64)$$

$$\begin{aligned} \int_{\mathbb{R}} p(t)q(t)dt &= \frac{1}{2\pi\sigma_P\sigma_Q} \int_{\mathbb{R}} e^{-\frac{(t-\mu_P)^2}{2\sigma_P^2} - \frac{(t-\mu_Q)^2}{2\sigma_Q^2}} dt = \frac{1}{2\pi\sigma_P\sigma_Q} \sqrt{2\pi} \frac{\sigma_P\sigma_Q}{\sqrt{\sigma_P^2 + \sigma_Q^2}} e^{-\frac{1}{2}\frac{(\mu_P - \mu_Q)^2}{\sigma_P^2 + \sigma_Q^2}} \\ &= \frac{1}{\sqrt{2\pi(\sigma_P^2 + \sigma_Q^2)}} e^{-\frac{(\mu_P - \mu_Q)^2}{2(\sigma_P^2 + \sigma_Q^2)}}. \end{aligned}$$

We can notice that this is the density of a Gaussian distribution with mean  $\mu_P$  and variance  $\sigma_P^2 + \sigma_Q^2$  evaluated at  $\mu_Q$ . We will also use a well-known notation for this value:  $\mathcal{N}(\mu_P | \mu_Q, \sigma_P^2 + \sigma_Q^2)$ . Here

$$\mathcal{N}(a | b, s^2) = \frac{1}{\sqrt{2\pi s^2}} e^{-\frac{(a-b)^2}{2s^2}}, \quad (65)$$

and

$$\text{QS}(P, Q) = -2\mathcal{N}(\mu_P | \mu_Q, \sigma_P^2 + \sigma_Q^2) + \frac{1}{2\sqrt{\pi}\sigma_P}. \quad (66)$$

Additionally, we present another useful identity involving the expected quadratic score for a Gaussian prediction and Gaussian outcome:

$$\int_{\mathbb{R}} p(t) q(t) dt = -\frac{\text{QS}(P, Q) + H(P)}{2}. \quad (67)$$

### B.3.4 QUADRATIC SCORE DIVERGENCE FOR TWO GAUSSIANS

$$d(P, Q) = \mathbb{E}_{Y \sim Q}[\text{QS}(P, Y)] - H(Q) = -2\mathbb{E}_{Y \sim Q}[p(Y)] - H(P) - H(Q). \quad (68)$$

We can already compute these terms using the results from the previous sections. We have the following:

$$-2\mathbb{E}_{Y \sim Q}[p(Y)] = -2\mathcal{N}(\mu_P | \mu_Q, \sigma_P^2 + \sigma_Q^2) = -\frac{2}{\sqrt{2\pi(\sigma_P^2 + \sigma_Q^2)}} e^{-\frac{(\mu_P - \mu_Q)^2}{2(\sigma_P^2 + \sigma_Q^2)}}. \quad (69)$$

$$H(P) = -\frac{1}{2\sqrt{\pi}\sigma_P}, \quad H(Q) = -\frac{1}{2\sqrt{\pi}\sigma_Q}. \quad (70)$$

Thus, we can express the quadratic score divergence for two Gaussians as follows:

$$d(P, Q) = -2\mathcal{N}(\mu_P | \mu_Q, \sigma_P^2 + \sigma_Q^2) + \frac{1}{2\sqrt{\pi}\sigma_P} + \frac{1}{2\sqrt{\pi}\sigma_Q}. \quad (71)$$

Another convenient expression for the quadratic score divergence is:

$$d(P, Q) = \mathbb{E}_{Y \sim Q}[\text{QS}(P, Y)] - H(Q) = -2 \int_{\mathbb{R}} p(t) q(t) dt + \int_{\mathbb{R}} p(t)^2 dt + \int_{\mathbb{R}} q(t)^2 dt = \int_{\mathbb{R}} (p(t) - q(t))^2 dt.$$

This divergence is symmetric:  $d(P, Q) = d(Q, P)$ . It is similar to CRPS score divergence, but this time we use densities instead of cumulative distribution functions.

### B.3.5 ENTROPY OF QUADRATIC SCORE FOR A GAUSSIAN MIXTURE

$$H(\hat{P}_{\text{ens}}) = -\int_{\mathbb{R}} \left( \frac{1}{M} \sum_{i=1}^M \hat{p}_i(y) \right)^2 dy = -\frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \mathcal{N}(\mu_i | \mu_j, \sigma_i^2 + \sigma_j^2). \quad (72)$$

Alternatively, we can express the entropy of quadratic score for a Gaussian mixture as follows:

$$H(\hat{P}_{\text{ens}}) = \frac{1}{M^2} \left[ \sum_{i=1}^M H(\hat{P}_i) - 2 \sum_{j < i} \mathcal{N}(\mu_i | \mu_j, \sigma_i^2 + \sigma_j^2) \right] \quad (73)$$

$$= -\frac{1}{M^2} \left[ \frac{1}{2\sqrt{\pi}} \sum_{i=1}^M \frac{1}{\sigma_i} + 2 \sum_{j < i} \mathcal{N}(\mu_i | \mu_j, \sigma_i^2 + \sigma_j^2) \right]. \quad (74)$$

Next we express the entropy of quadratic score for a Gaussian mixture in terms of the expected quadratic score and entropy of its components:

$$\begin{aligned} H(\hat{P}_{\text{ens}}) &= -\frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \int_{\mathbb{R}} \hat{p}_i(y) \hat{p}_j(y) dy = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \frac{\text{QS}(\hat{P}_i, \hat{P}_j) + H(\hat{P}_i)}{2} \\ &= \frac{1}{2M} \sum_{i=1}^M H(\hat{P}_i) + \frac{1}{2M^2} \sum_{i=1}^M \sum_{j=1}^M \text{QS}(\hat{P}_i, \hat{P}_j) = \frac{M+1}{2M^2} \sum_{i=1}^M H(\hat{P}_i) + \frac{1}{2M^2} \sum_{i \neq j} \text{QS}(\hat{P}_i, \hat{P}_j). \end{aligned}$$

### B.3.6 QUADRATIC SCORE DIVERGENCE BETWEEN A GAUSSIAN AND A GAUSSIAN MIXTURE

$$\begin{aligned}
d(\widehat{P}_i, \widehat{P}_{\text{ens}}) &= \mathbb{E}_{Y \sim \widehat{P}_{\text{ens}}} [\text{QS}(\widehat{P}_i, Y)] - H(\widehat{P}_{\text{ens}}) = \frac{1}{M} \sum_{j=1}^M \mathbb{E}_{Y \sim \widehat{P}_j} [\text{QS}(\widehat{P}_i, Y)] - H(\widehat{P}_{\text{ens}}) \\
&= \frac{1}{M} \sum_{j=1}^M \text{QS}(\widehat{P}_i, \widehat{P}_j) - H(\widehat{P}_{\text{ens}}) = -\frac{2}{M} \sum_{j=1}^M \mathcal{N}(\mu_i | \mu_j, \sigma_i^2 + \sigma_j^2) - H(\widehat{P}_i) - H(\widehat{P}_{\text{ens}}).
\end{aligned} \tag{75}$$

## C PROPER SCORING RULES FROM CONSISTENT SCORING FUNCTIONS

A closely related concept to proper scoring rule is the *consistent scoring function*. In some situations a probabilistic forecast is not possible or not needed (e.g., a point prediction is required by some regulation); for this setup a similar theory was developed (Savage, 1971; Gneiting, 2009). There is one important difference when evaluating point forecasts: since we still follow the random data model (the true outcome  $y$  is stochastic, comes from distribution  $P$ ), we need to select a functional  $T(P)$  of the true data distribution that we want to estimate. You can refer to (Gneiting, 2009) for a detailed discussion, and in the next subsection we provide a brief overview of the main definitions and results.

### C.1 INTRODUCTION

**Definition 3** (Definition 2.1 from (Gneiting, 2009)). Scoring function  $S$  is consistent for the functional  $T$  relative to the class of distributions  $\mathcal{P}$  if

$$\mathbb{E}_{Y \sim P} [S(t, Y)] \leq \mathbb{E}_{Y \sim P} [S(x, Y)], \tag{76}$$

for all distributions  $P \in \mathcal{P}$ , all  $t \in T(P)$  and all  $x \in \text{dom}(Y)$ .

One of the most famous results in this area is the following.

**Theorem 4** (Savage (Savage, 1971), 1971). *Every scoring function that is consistent for the mean functional  $T(P) = \mathbb{E}_P[Y]$  admits a representation as a Bregman divergence:*

$$S(x, y) = \varphi(y) - \varphi(x) - \varphi'(x)(y - x), \tag{77}$$

where  $\varphi$  is a convex function, and  $\varphi'$  is its subgradient. Such  $S(x, y)$  are also called Bregman functions.

For example, the quadratic score  $S(x, y) = (y - x)^2$  is a Bregman function with  $\varphi(t) = t^2$ .

**Theorem 5** (Theorem 2.2 from (Gneiting, 2009)). *Score  $S$  is consistent for  $T$  if and only if any  $t \in T(P)$  is an optimal point prediction:  $t \in y_{\text{Bayes}}$ .*

Optimal point prediction in this setting is the Bayes act  $\widehat{y}_{\text{Bayes}}$  and can be expressed as follows:

$$\widehat{y}_{\text{Bayes}} = \arg \min_x \mathbb{E}_{Y \sim P} [S(x, Y)]. \tag{78}$$

In the case of the mean functional,  $y_{\text{Bayes}} = T(P) = \mathbb{E}_{Y \sim P} [Y]$  which also follows from the representation of the consistent score as a Bregman function. Now we are ready to introduce the main tool. If we have a scoring function that is consistent for some functional, then we can use it to construct a (rather simple) proper scoring rule:

**Theorem 6** (Theorem 2.3 from (Gneiting, 2009)). *If  $S_T$  is consistent for the functional  $T(\cdot)$ , then the following scoring rule is proper:*

$$S(\widehat{P}, y) = S_T(T(\widehat{P}), y).$$

## 1512 C.2 SE SCORE FROM QUADRATIC SCORING FUNCTION

1513  
1514 Based on the quadratic scoring function, which is consistent for the mean, we introduce the following  
1515 proper scoring rule for continuous distributions:

$$1516 \text{SE}(\widehat{P}, y) = (y - \mathbb{E}_{Y \sim \widehat{P}}[Y])^2,$$

1517  
1518 and call it the SE score. This score is proper, but not strictly proper, since any distribution with the  
1519 same mean will have the same score.

## 1521 C.3 RISK COMPONENTS FOR SE SCORE

1522  
1523 The risk components for this score are as follows:

- 1524 • **Total risk:**  $R_{\text{Tot}}(\widehat{P}, P) = \mathbb{E}_{Y \sim P}[\text{SE}(\widehat{P}, Y)] = \mathbb{E}_{Y \sim P}[(Y - \mathbb{E}_{X \sim \widehat{P}}[X])^2]$ , where  
1525  
1526  $\mathbb{E}_{Y \sim P}[(Y - \mathbb{E}_{X \sim \widehat{P}}[X])^2] = \mathbb{E}_{Y \sim P}[Y^2] - 2\mathbb{E}_{Y \sim P}[Y] \cdot \mathbb{E}_{X \sim \widehat{P}}[X] + (\mathbb{E}_{X \sim \widehat{P}}[X])^2$   
1527  
1528  $= (\mathbb{E}_{X \sim \widehat{P}}[X] - \mathbb{E}_{Y \sim P}[Y])^2 + \mathbb{E}_{Y \sim P}[Y^2] - (\mathbb{E}_{Y \sim P}[Y])^2$   
1529  
1530  $= (\mathbb{E}_{X \sim \widehat{P}}[X] - \mathbb{E}_{Y \sim P}[Y])^2 + \text{Var}_{Y \sim P}[Y]$ .
- 1531  $R_{\text{Tot}}(\widehat{P}, P) = (\mathbb{E}_{X \sim \widehat{P}}[X] - \mathbb{E}_{Y \sim P}[Y])^2 + \text{Var}_{Y \sim P}[Y]$ .
- 1532 • **Bayes risk:**  $R_{\text{Bayes}}(P) = H_{\text{SE}}(P) = \mathbb{E}_{Y \sim P}[\text{SE}(P, Y)] = \text{Var}_{Y \sim P}[Y]$ .
- 1533 • **Excess risk:**  $R_{\text{Exc}}(\widehat{P}, P) = d_{\text{SE}}(\widehat{P}, P) = R_{\text{Tot}}(\widehat{P}, P) - R_{\text{Bayes}}(P) =$   
1534  
1535  $(\mathbb{E}_{X \sim \widehat{P}}[X] - \mathbb{E}_{Y \sim P}[Y])^2$ .

1536  
1537 In other words, the SE score corresponds to the following entropy and divergence functions:

$$1538 H_{\text{SE}}(P) = \text{Var}_{Y \sim P}[Y], \quad d_{\text{SE}}(\widehat{P}, P) = (\mathbb{E}_{X \sim \widehat{P}}[X] - \mathbb{E}_{Y \sim P}[Y])^2.$$

## 1541 D CHARACTERIZATION OF PROPER SCORING RULES AND GENERAL FORM OF

### 1542 RISK ESTIMATES

1543  
1544 The following theorem provides a general form of a proper scoring rule using the notion of a convex  
1545 function on the space of probability measures:

1546 **Definition 7.** A scoring rule  $S$  is *regular* if  $H(P) = S(P, P)$  is finite and  $S(P, Q) > -\infty, \forall P, Q \in \mathcal{P}$ .

1547 **Theorem 8** (Theorem 1 from (Gneiting & Raftery, 2007), notation from Theorem 12 in (Waghmare  
1548 & Ziegel, 2025)). *A regular scoring rule is proper if and only if there exists a concave function*  
1549  $H: \mathcal{P} \rightarrow \mathbb{R}$  *such that*

$$1550 S(P, y) = H(P) + \langle h_P, \delta_y - P \rangle = H(P) + h_P(y) - \int h_P(t) dP(t) \quad (79)$$

1551  
1552 for every  $P \in \mathcal{P}$  and  $y \in \mathcal{Y}$ , where  $h_P$  is a supergradient of  $H$  at  $P$  and  $\langle h_P, Q \rangle = \int h_P(t) dQ(t)$ .

1553  
1554 From this theorem we can obtain the following general expression of the divergence for a proper  
1555 scoring rule:

$$1556 d(P, Q) = H(P) - H(Q) - \langle h_P, P - Q \rangle. \quad (80)$$

1557  
1558 The concave function of Theorem 8 coincides with the corresponding score entropy  $H$ . In case of our  
1559 ensemble setting, this leads to the following relation between the risk estimates:

$$1560 \widehat{R}_{\text{Bayes}}^1 = \frac{1}{M} \sum_{i=1}^M H(\widehat{P}_i) \leq H \left( \frac{1}{M} \sum_{i=1}^M \widehat{P}_i \right) \quad (81)$$

$$1561 = H(\widehat{P}_{\text{ens}}) = \widehat{R}_{\text{Bayes}}^2$$

Table 8: Representations of some proper scoring rules

Score	Entropy $H(P)$	Supergradient map $h_P(y)$
CRPS	$\frac{1}{2}\mathbb{E}_{X, X' \sim P}  X - X' $	$\mathbb{E}_{X \sim P}  X - y $
LS	$-\int_{\mathbb{R}} p(t) \log p(t) dt$	$-\log p(y)$
QS	$-\int_{\mathbb{R}} p(t)^2 dt$	$-2p(y)$
SE	$\text{Var}_{Y \sim P} [Y]$	$y^2 - 2y\mathbb{E}_{X \sim P} [X]$

For the Logarithmic score, we know that:

$$\text{LS: } \widehat{R}_{\text{Bayes}}^3 = H(\widehat{P}_*) \geq H(\widehat{P}_{\text{ens}}) = \widehat{R}_{\text{Bayes}}^2, \quad (82)$$

since among the distributions with the same variance (in this case, both have variance  $\sigma_*^2$ ), the Gaussian distribution has the highest (Shannon) entropy. Combining these results, we get:

$$\text{LS: } \widehat{R}_{\text{Bayes}}^1 \leq \widehat{R}_{\text{Bayes}}^2 \leq \widehat{R}_{\text{Bayes}}^3. \quad (83)$$

Whether this result holds for other scores or in general remains open. We think that it depends on the relationship between the particular score’s entropy and the Gaussian distribution that we used for our ensemble estimate.

#### CRPS.

$$\widehat{R}_{\text{Bayes}}^3 = \sqrt{\frac{\sigma_*^2}{\pi}} = \frac{1}{\sqrt{\pi}} \sqrt{\frac{1}{M} \sum_{i=1}^M \sigma_i^2 + \text{Var}[\mu_i]} \geq \frac{1}{\sqrt{\pi}} \sqrt{\frac{1}{M} \sum_{i=1}^M \sigma_i^2} \geq \frac{1}{\sqrt{\pi}} \frac{1}{M} \sum_{i=1}^M \sigma_i = \widehat{R}_{\text{Bayes}}^1.$$

#### QS.

$$\begin{aligned} \widehat{R}_{\text{Bayes}}^3 &= -\frac{1}{2\sqrt{\pi}\sigma_*} = -\frac{1}{2\sqrt{\pi}} \frac{1}{\sqrt{\frac{1}{M} \sum_{i=1}^M \sigma_i^2 + \text{Var}[\mu_i]}} \geq \\ &= -\frac{1}{2\sqrt{\pi}} \frac{1}{\sqrt{\frac{1}{M} \sum_{i=1}^M \sigma_i^2}} \geq -\frac{1}{2\sqrt{\pi}} \frac{1}{M} \sum_{i=1}^M \frac{1}{\sigma_i} = \widehat{R}_{\text{Bayes}}^1. \end{aligned}$$

Here we have applied Jensen’s inequality to the function  $f(x) = -\frac{1}{\sqrt{x}}$ ,  $x_i = \sigma_i^2$ .

#### SE.

$$\widehat{R}_{\text{Bayes}}^3 = \sigma_*^2 = \frac{1}{M} \sum_{i=1}^M \sigma_i^2 + \text{Var}[\mu_i] \geq \frac{1}{M} \sum_{i=1}^M \sigma_i^2 = \widehat{R}_{\text{Bayes}}^1.$$

#### D.1 ADDITIONAL RESULTS FOR EXCESS RISK

$$\widehat{R}_{\text{Exc}}^{2,1} = \frac{1}{M} \sum_{i=1}^M d(\widehat{P}_{\text{ens}}, \widehat{P}_i) = \frac{1}{M} \sum_{i=1}^M \left[ H(\widehat{P}_{\text{ens}}) - H(\widehat{P}_i) - \langle h_{\widehat{P}_{\text{ens}}}, \widehat{P}_{\text{ens}} - \widehat{P}_i \rangle \right] \quad (84)$$

$$= H(\widehat{P}_{\text{ens}}) - \frac{1}{M} \sum_{i=1}^M H(\widehat{P}_i) - \langle h_{\widehat{P}_{\text{ens}}}, \widehat{P}_{\text{ens}} - \frac{1}{M} \sum_{i=1}^M \widehat{P}_i \rangle = \widehat{R}_{\text{Bayes}}^2 - \widehat{R}_{\text{Bayes}}^1 \geq 0. \quad (85)$$

1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673

$$\begin{aligned}
 \widehat{R}_{\text{Exc}}^{1,1} &= \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M d(\widehat{P}_i, \widehat{P}_j) = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \left[ H(\widehat{P}_i) - H(\widehat{P}_j) - \langle h_{\widehat{P}_i}, \widehat{P}_i - \widehat{P}_j \rangle \right] \\
 &= -\frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \langle h_{\widehat{P}_i}, \widehat{P}_i - \widehat{P}_j \rangle = -\frac{1}{M} \sum_{i=1}^M \langle h_{\widehat{P}_i}, \widehat{P}_i - \frac{1}{M} \sum_{j=1}^M \widehat{P}_j \rangle = \frac{1}{M} \sum_{i=1}^M \langle h_{\widehat{P}_i}, \widehat{P}_{\text{ens}} - \widehat{P}_i \rangle \\
 &\geq \frac{1}{M} \sum_{i=1}^M \left[ H(\widehat{P}_{\text{ens}}) - H(\widehat{P}_i) \right] = H(\widehat{P}_{\text{ens}}) - \frac{1}{M} \sum_{i=1}^M H(\widehat{P}_i) = \widehat{R}_{\text{Bayes}}^2 - \widehat{R}_{\text{Bayes}}^1 = \widehat{R}_{\text{Exc}}^{2,1}.
 \end{aligned} \tag{86}$$

$$\begin{aligned}
 \widehat{R}_{\text{Exc}}^{1,2} &= \frac{1}{M} \sum_{i=1}^M d(\widehat{P}_i, \widehat{P}_{\text{ens}}) = \frac{1}{M} \sum_{i=1}^M \left[ H(\widehat{P}_i) - H(\widehat{P}_{\text{ens}}) - \langle h_{\widehat{P}_i}, \widehat{P}_i - \widehat{P}_{\text{ens}} \rangle \right] \\
 &= \frac{1}{M} \sum_{i=1}^M H(\widehat{P}_i) - H(\widehat{P}_{\text{ens}}) + \frac{1}{M} \sum_{i=1}^M \langle h_{\widehat{P}_i}, \widehat{P}_{\text{ens}} - \widehat{P}_i \rangle = \widehat{R}_{\text{Exc}}^{1,1} - \widehat{R}_{\text{Exc}}^{2,1} \\
 &\implies \widehat{R}_{\text{Exc}}^{1,1} = \widehat{R}_{\text{Exc}}^{2,1} + \widehat{R}_{\text{Exc}}^{1,2}.
 \end{aligned} \tag{87}$$

## E EXPERIMENTS

The code for reproducing all the experiments in this paper is provided at TBD

### E.1 TRAINING OBJECTIVE

In our experiments, we consider neural networks predicting the mean and variance of an output distribution  $\mathcal{N}(\mu(x; w), \sigma^2(x; w))$  over  $y$  for a given input  $x$ . The standard way of optimizing a model with parameters  $w$  to predict a faithful output distribution is by minimizing the log-likelihood of the output distribution on the training set  $D$  (Nix & Weigend, 1994; Lakshminarayanan et al., 2017). Formulated as negative log-likelihood, the element-wise loss is given by

$$\mathcal{L}(w) = \frac{1}{2} \log \sigma^2(x; w) + \frac{(y - \mu(x; w))^2}{2\sigma^2(x; w)} + \text{const.} \quad (88)$$

However, training both mean and variance networks in conjunction is known to suffer from instability. Therefore, there have been multiple attempts to stabilize training and obtain both accurate mean predictions and calibrated variance (Skafte et al., 2019; Seitzer et al., 2022; Immer et al., 2023). We follow the approach of Immer et al. (2023) in using the natural parameterization of the Gaussian to reformulate the optimization problem and obtain more stable optimization properties. Instead of predicting mean  $\mu$  and variance  $\sigma^2$ , in the natural parameterization  $\eta_1 = \frac{\mu}{\sigma^2}$  and  $\eta_2 = -\frac{1}{2\sigma^2}$  are predicted, where one needs to make sure that  $\eta_2 < 0$ . Then, the negative log-likelihood is given by

$$\mathcal{L}(w) = - \begin{pmatrix} \eta_1(x; w) \\ \eta_2(x; w) \end{pmatrix}^T \begin{pmatrix} y \\ y^2 \end{pmatrix} - \frac{\eta_1(x; w)^2}{4\eta_2(x; w)} - \frac{1}{2} \log(-2\eta_2(x; w)) + \text{const.} \quad (89)$$

We compare the two losses on a simple regression dataset using a three layer neural networks consisting of feed-forward layers with hidden dimension 32 and ReLU activation. Data points are drawn according to a sine function with larger noise in the right open half-plane. While we observe in Figure 4 that using Equation (89) needs more gradient steps than optimization with Equation (88), training and validation losses are much more stable. Furthermore, the resulting behavior is arguable more favorable for the natural parameterization (see Figure 5). For the standard parameterization, the predicted variance collapses away from data, while it stays at least relatively constant for the natural parameterization. Additionally, for the standard parameterization the predicted mean follows a linear trend, while for the natural parameterization it resorts to a constant value. As the natural parameterization lead to much more stable results throughout all our experiments, we exclusively consider the natural parameterization as training objective in the reported results.

### E.2 SYNTHETIC EXPERIMENT

The ground-truth conditional is a heteroscedastic two-component mixture,

$$p^*(y | x) = \pi(x)P_1(y | x) + (1 - \pi(x))P_2(y | x),$$

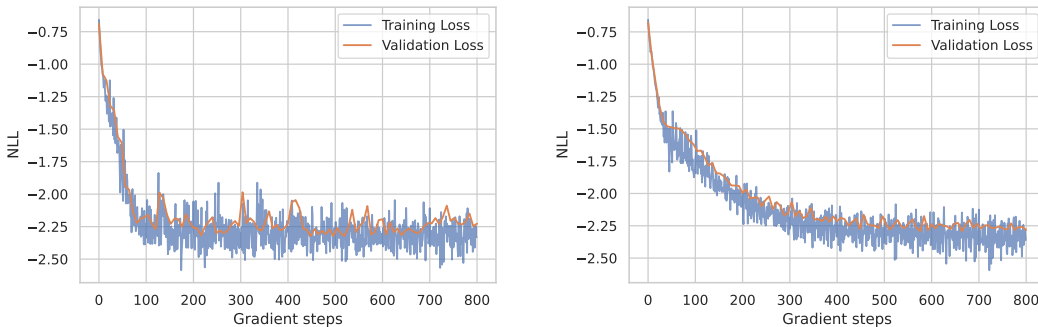


Figure 4: Training and validation loss when minimizing the Gaussian NLL on a synthetic regression example. **Left:** Training with the standard parameterization (equation (88)). **Right:** Training with the natural parameterization (equation (89)).

1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

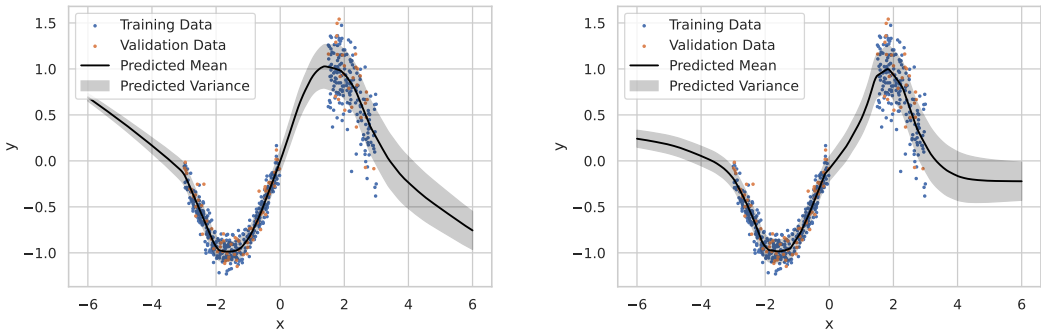


Figure 5: Best model when minimizing the Gaussian NLL on a synthetic regression example. **Left:** Predicting with the standard parameterization (equation (88)). **Right:** Predicting with the natural parameterization (equation (89)).

with mixing weight and component means

$$\begin{aligned} \pi(x) &= \frac{1}{1 + \exp(1.2x)}, \\ \mu_1(x) &= \frac{x}{3} + 1.2 \sin(0.8x), \\ \mu_2(x) &= \frac{x}{3} - 1.2 \cos(0.8x), \end{aligned}$$

and input-dependent noise scale

$$\sigma(x) = 0.12 + 0.28(0.5 + 0.5 \sin(0.7x))^2.$$

We take the components to be Gaussian with shared heteroscedastic variance,

$$P_1(y | x) = \mathcal{N}(y; \mu_1(x), \sigma^2(x)), \quad P_2(y | x) = \mathcal{N}(y; \mu_2(x), \sigma^2(x)),$$

equivalently  $y = \mu_k(x) + \epsilon\sigma(x)$  with  $\epsilon \sim \mathcal{N}(0, 1)$  and  $k \in \{1, 2\}$  drawn according to  $\pi(x)$ . Figure 6 shows the ground-truth conditional and samples.

We draw  $n = 1200$  training pairs from  $p^*(y | x)$  and fit a regression network. The backbone is a fully connected MLP with two hidden layers of width 8 and SiLU (Hendrycks & Gimpel, 2016). On top of the backbone, we use two output heads that parameterize a Gaussian predictive distribution via its natural parameters (Immer et al., 2023), and we train by maximizing the corresponding Gaussian log-likelihood. We form an ensemble of 10 independently initialized models, each trained for 100 epochs with Adam (Kingma & Ba, 2014).

Given the trained ensemble, we compute our uncertainty scores. Figure 2 in the main part of the paper reports one instance of our framework using the logarithmic score: each dot shows the ensemble-averaged predictive mean, and the color intensity encodes the corresponding uncertainty (see color bar).

Below we extend the main-text visualization to additional scoring rules (CRPS, SE, and Quadratic). Across all rules, Bayes risk is highest in regions where the mapping is least determined, i.e., where the two generating curves are sampled with comparable probability, while excess risk rises for inputs outside the training data support, reflecting increased epistemic uncertainty. Although absolute magnitudes vary across rules (a score-specific effect), the spatial patterns are consistent. Extended plots are shown in Figures 7-8.

### E.3 SELECTIVE PREDICTION

**Datasets.** We evaluate on seven datasets. First, we introduce `dots`, where inputs are  $32 \times 32$  grayscale images containing non-overlapping dots and targets are the corresponding dot counts. Second, we introduce `arrow`, where inputs are  $32 \times 32$  grayscale images of left-pointing arrows and targets are the arrow angles. Third, we derive five datasets from Cityscapes (Cordts et al.,

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

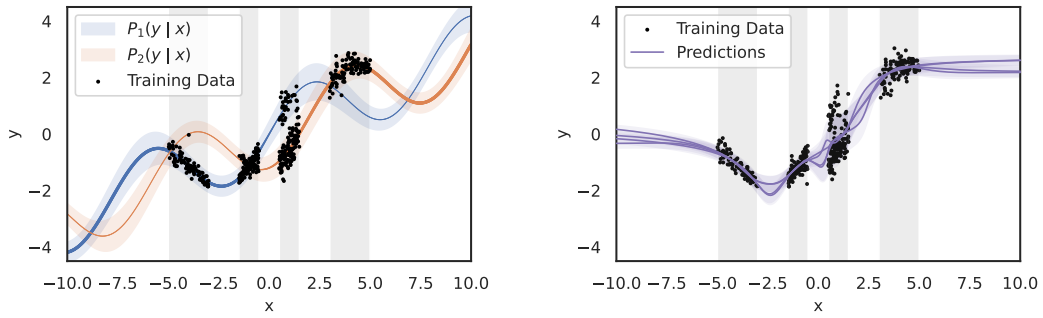


Figure 6: Synthetic regression setup. **Left:** Mixture components  $P_1$  and  $P_2$  of the ground-truth conditional  $p^*(y|x)$ ; curve thickness encodes the mixing weights  $\pi(x)$  and  $1 - \pi(x)$ , respectively. **Right:** Training samples drawn from  $p^*$  and predictive distributions of individual ensemble members.

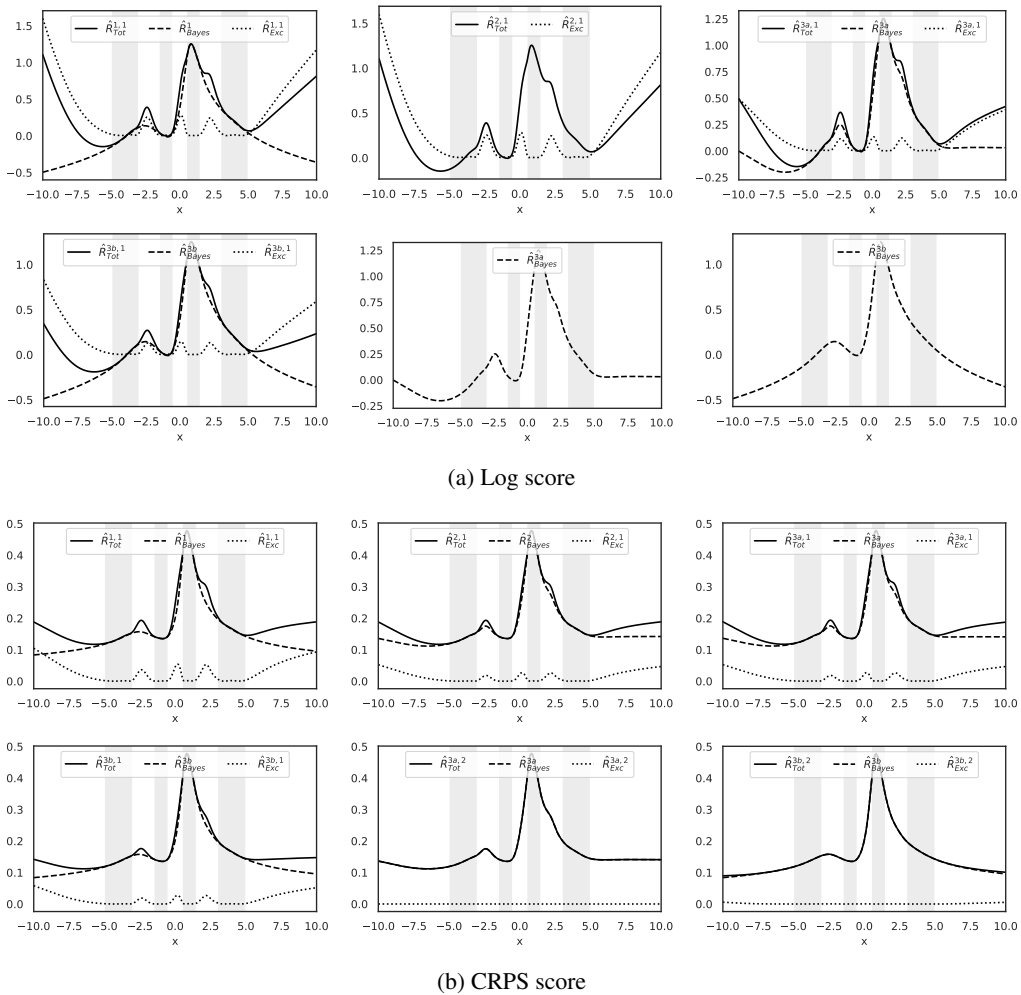


Figure 7: Synthetic 1D experiment (part 1): Log and CRPS score.

2016) by tiling each image into  $224 \times 224$  RGB crops and using the provided segmentation maps to define regression targets: the number of pixels belonging to a given class. The classes are street, building, sky, car, and vegetation.

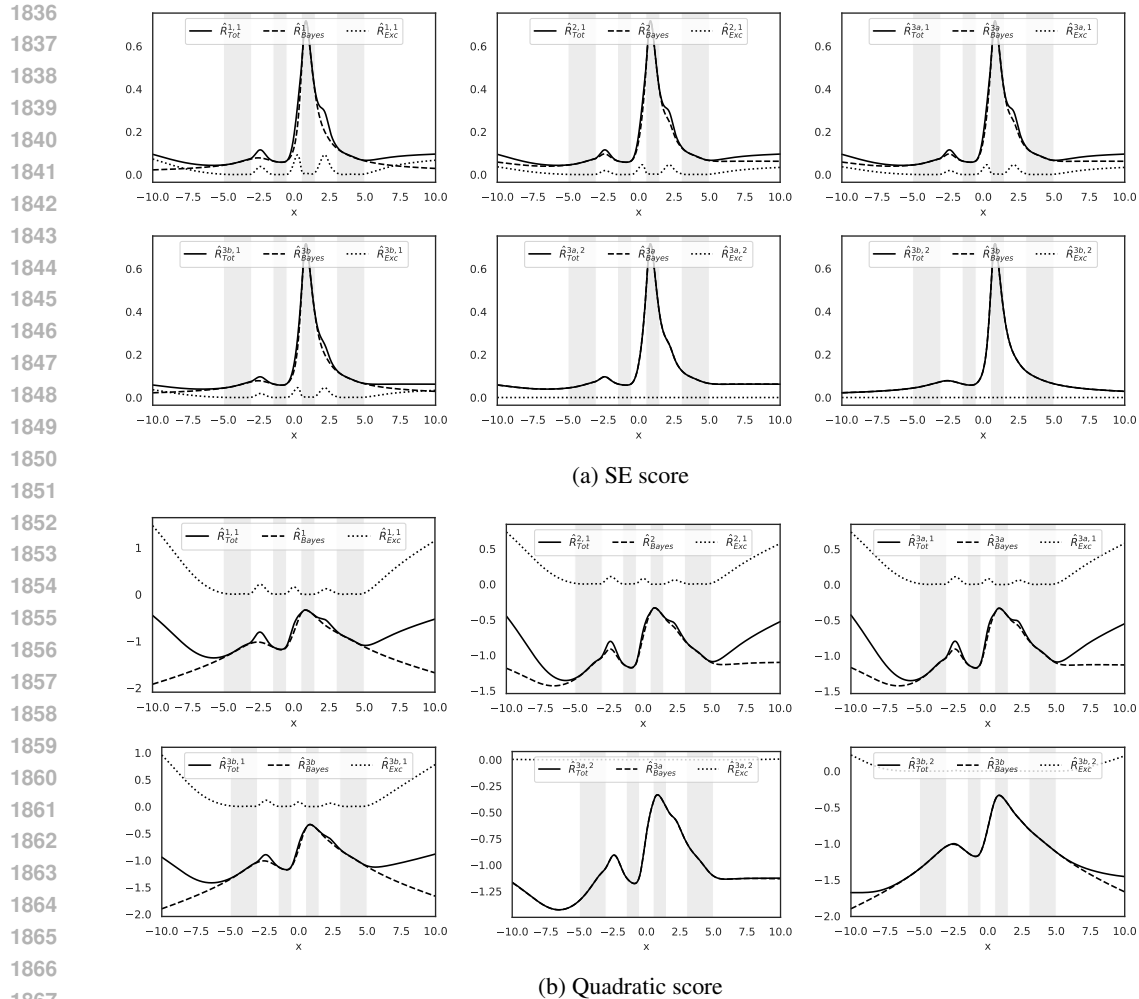


Figure 8: Synthetic 1D experiment (part 2): SE and Quadratic score.

For train–test splits, we generate new samples for `dots` and `arrow`, and use the official Cityscapes test split (Cordts et al., 2016) for the remaining datasets.

**Training details.** For `dots` and `arrow`, we use a small CNN with three convolutional layers followed by two fully connected layers and ReLU activations. Models are trained with Adam (Kingma & Ba, 2014), learning rate  $10^{-3}$ , and batch size 64. For the Cityscapes-derived datasets, we train ResNet-18 (He et al., 2016) with the same optimizer settings. Unless stated otherwise, we report results from ensembles of 10 independently trained models per dataset.

**Extended results.** Selective prediction is quantified using the prediction–reject ratio (PRR; Malinin & Gales, 2021). Although introduced for classification, PRR extends directly to regression by fixing the performance metric used to compute areas; we use mean squared error (MSE). We consider retention rates (i.e.,  $1 - \text{rejection}$ ) from 0.5 to 1. Example retention curves with PRR are shown in Figure 9 (insets), and full per-dataset results appear in Table 9.

#### E.4 OUT-OF-DISTRIBUTION DETECTION

**Datasets.** We construct inputs as  $64 \times 64$  grayscale mosaics of four MNIST digits (Lecun et al., 1998) arranged (upper-left, upper-right, bottom-left, bottom-right) to form a four-digit number; the target is this number. For out-of-distribution (OOD) data, we substitute each quadrant with images

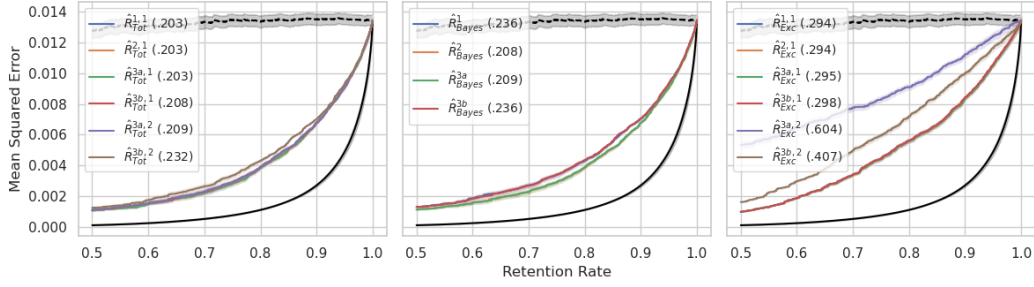


Figure 9: Exemplary retention curve for `street` dataset, PRRs reported in the inset. Black solid line is the optimal, oracle retention curve sorted by the actual MSE of each prediction. The black dashed line is the random baseline.

Table 9: Results for selective prediction (PRR $\downarrow$ ) on different datasets  $\mathcal{D}$  for different scoring rules (SR).

$\mathcal{D}$	SR	$\hat{R}_{Tot}^{1,1}$	$\hat{R}_{Tot}^{2,1}$	$\hat{R}_{3a,1}^{3a,1}$	$\hat{R}_{3b,1}^{3b,1}$	$\hat{R}_{3a,2}^{3a,2}$	$\hat{R}_{3b,2}^{3b,2}$	$\hat{R}_{Bayes}^1$	$\hat{R}_{Bayes}^2$	$\hat{R}_{Bayes}^{3a}$	$\hat{R}_{Bayes}^{3b}$	$\hat{R}_{Exc}^{1,1}$	$\hat{R}_{Exc}^{2,1}$	$\hat{R}_{Exc}^{3a,1}$	$\hat{R}_{Exc}^{3b,1}$	$\hat{R}_{Exc}^{3a,2}$	$\hat{R}_{Exc}^{3b,2}$
dots	CRPS	0.631 ±0.002	0.631 ±0.002	0.631 ±0.002	0.644 ±0.002	0.644 ±0.002	0.664 ±0.002	0.664 ±0.002	0.644 ±0.002	0.644 ±0.002	0.664 ±0.002	0.704 ±0.001	0.704 ±0.001	0.704 ±0.001	0.705 ±0.001	0.782 ±0.002	0.716 ±0.001
	Log	0.631 ±0.002	0.631 ±0.002	0.631 ±0.002	0.643 ±0.002	-	-	0.664 ±0.002	-	0.644 ±0.002	0.664 ±0.002	0.729 ±0.001	0.729 ±0.001	0.729 ±0.001	0.729 ±0.001	-	-
	SE	0.632 ±0.002	0.632 ±0.002	0.632 ±0.002	0.644 ±0.002	0.644 ±0.002	0.664 ±0.002	0.664 ±0.002	0.644 ±0.002	0.644 ±0.002	0.664 ±0.002	0.685 ±0.001	0.685 ±0.001	0.685 ±0.001	0.685 ±0.001	-	-
	Quad.	0.631 ±0.002	0.631 ±0.002	0.631 ±0.002	0.643 ±0.002	0.644 ±0.002	0.663 ±0.002	0.664 ±0.002	0.644 ±0.002	0.644 ±0.002	0.664 ±0.002	0.760 ±0.002	0.760 ±0.002	0.760 ±0.002	0.759 ±0.002	0.805 ±0.002	0.743 ±0.001
arrow	CRPS	0.776 ±0.017	0.776 ±0.017	0.776 ±0.017	0.785 ±0.018	0.786 ±0.018	0.806 ±0.018	0.810 ±0.018	0.786 ±0.018	0.786 ±0.018	0.810 ±0.017	0.756 ±0.012	0.756 ±0.012	0.756 ±0.012	0.756 ±0.012	0.804 ±0.014	0.762 ±0.011
	Log	0.773 ±0.017	0.773 ±0.017	0.775 ±0.017	0.783 ±0.018	-	-	0.811 ±0.017	-	0.786 ±0.018	0.810 ±0.017	0.765 ±0.010	0.765 ±0.010	0.765 ±0.010	0.766 ±0.011	-	-
	SE	0.777 ±0.017	0.777 ±0.017	0.777 ±0.017	0.786 ±0.018	0.786 ±0.018	0.810 ±0.017	0.810 ±0.018	0.786 ±0.018	0.786 ±0.018	0.810 ±0.017	0.753 ±0.013	0.753 ±0.013	0.753 ±0.013	0.753 ±0.013	-	-
	Quad.	0.774 ±0.017	0.774 ±0.017	0.774 ±0.017	0.784 ±0.018	0.785 ±0.018	0.801 ±0.018	0.811 ±0.017	0.786 ±0.018	0.786 ±0.018	0.810 ±0.017	0.783 ±0.010	0.783 ±0.010	0.783 ±0.010	0.782 ±0.010	0.815 ±0.014	0.776 ±0.010
street	CRPS	0.203 ±0.013	0.203 ±0.013	0.204 ±0.013	0.209 ±0.014	0.209 ±0.014	0.232 ±0.014	0.236 ±0.014	0.209 ±0.014	0.209 ±0.014	0.237 ±0.014	0.295 ±0.012	0.295 ±0.012	0.295 ±0.012	0.298 ±0.012	0.605 ±0.016	0.407 ±0.012
	Log	0.215 ±0.013	0.215 ±0.013	0.206 ±0.013	0.209 ±0.014	-	-	0.237 ±0.014	-	0.209 ±0.014	0.237 ±0.014	0.479 ±0.014	0.479 ±0.014	0.488 ±0.014	0.488 ±0.014	-	-
	SE	0.204 ±0.013	0.204 ±0.013	0.204 ±0.013	0.209 ±0.014	0.209 ±0.014	0.237 ±0.014	0.237 ±0.014	0.209 ±0.014	0.209 ±0.014	0.237 ±0.014	0.238 ±0.011	0.238 ±0.011	0.238 ±0.011	0.238 ±0.011	-	-
	Quad.	0.206 ±0.013	0.206 ±0.013	0.211 ±0.013	0.208 ±0.014	0.210 ±0.014	0.228 ±0.014	0.238 ±0.014	0.208 ±0.014	0.209 ±0.014	0.237 ±0.014	0.896 ±0.010	0.896 ±0.010	0.907 ±0.011	0.883 ±0.011	1.063 ±0.011	0.787 ±0.013
building	CRPS	0.227 ±0.012	0.227 ±0.012	0.227 ±0.012	0.235 ±0.012	0.236 ±0.012	0.264 ±0.012	0.273 ±0.012	0.237 ±0.012	0.236 ±0.012	0.270 ±0.012	0.240 ±0.011	0.240 ±0.011	0.242 ±0.011	0.242 ±0.011	0.551 ±0.012	0.318 ±0.013
	Log	0.227 ±0.012	0.227 ±0.012	0.227 ±0.012	0.232 ±0.012	-	-	0.277 ±0.012	-	0.236 ±0.012	0.270 ±0.012	0.336 ±0.013	0.336 ±0.013	0.338 ±0.013	0.332 ±0.013	-	-
	SE	0.228 ±0.012	0.228 ±0.012	0.228 ±0.012	0.236 ±0.012	0.236 ±0.012	0.270 ±0.012	0.270 ±0.012	0.236 ±0.012	0.236 ±0.012	0.270 ±0.012	0.229 ±0.011	0.229 ±0.011	0.229 ±0.011	0.229 ±0.011	-	-
	Quad.	0.226 ±0.012	0.226 ±0.012	0.227 ±0.012	0.233 ±0.012	0.235 ±0.012	0.257 ±0.012	0.282 ±0.011	0.242 ±0.012	0.236 ±0.012	0.270 ±0.012	0.517 ±0.012	0.517 ±0.012	0.538 ±0.012	0.523 ±0.012	0.888 ±0.011	0.568 ±0.013
sky	CRPS	0.043 ±0.003	0.043 ±0.003	0.044 ±0.003	0.045 ±0.004	0.045 ±0.004	0.050 ±0.006	0.051 ±0.006	0.044 ±0.004	0.045 ±0.004	0.052 ±0.006	0.044 ±0.001	0.044 ±0.001	0.045 ±0.001	0.044 ±0.001	0.099 ±0.004	0.050 ±0.001
	Log	0.046 ±0.003	0.046 ±0.003	0.045 ±0.003	0.046 ±0.004	-	-	0.050 ±0.006	-	0.045 ±0.004	0.052 ±0.006	0.054 ±0.001	0.054 ±0.001	0.057 ±0.001	0.056 ±0.001	-	-
	SE	0.044 ±0.004	0.044 ±0.004	0.044 ±0.004	0.045 ±0.004	0.045 ±0.004	0.052 ±0.006	0.052 ±0.006	0.045 ±0.004	0.045 ±0.004	0.052 ±0.006	0.043 ±0.002	0.043 ±0.002	0.043 ±0.002	0.043 ±0.002	-	-
	Quad.	0.041 ±0.003	0.041 ±0.003	0.045 ±0.003	0.044 ±0.004	0.046 ±0.004	0.048 ±0.006	0.049 ±0.006	0.042 ±0.004	0.045 ±0.004	0.052 ±0.006	0.075 ±0.003	0.075 ±0.003	0.079 ±0.003	0.073 ±0.003	0.172 ±0.004	0.072 ±0.002
car	CRPS	0.149 ±0.003	0.149 ±0.003	0.154 ±0.003	0.162 ±0.003	0.162 ±0.003	0.189 ±0.004	0.186 ±0.005	0.156 ±0.003	0.161 ±0.003	0.196 ±0.005	0.151 ±0.006	0.151 ±0.006	0.154 ±0.006	0.154 ±0.006	0.347 ±0.008	0.195 ±0.008
	Log	0.176 ±0.004	0.176 ±0.004	0.163 ±0.003	0.172 ±0.003	-	-	0.177 ±0.004	-	0.161 ±0.003	0.196 ±0.005	0.219 ±0.007	0.219 ±0.007	0.226 ±0.009	0.219 ±0.009	-	-
	SE	0.153 ±0.003	0.153 ±0.003	0.153 ±0.003	0.161 ±0.003	0.161 ±0.003	0.196 ±0.005	0.196 ±0.005	0.161 ±0.003	0.161 ±0.003	0.196 ±0.005	0.141 ±0.005	0.141 ±0.005	0.141 ±0.005	0.141 ±0.005	-	-
	Quad.	0.153 ±0.005	0.153 ±0.005	0.174 ±0.004	0.174 ±0.003	0.170 ±0.003	0.187 ±0.004	0.170 ±0.004	0.148 ±0.003	0.161 ±0.003	0.196 ±0.005	0.352 ±0.016	0.352 ±0.016	0.377 ±0.014	0.342 ±0.015	0.617 ±0.012	0.334 ±0.013
vegetation	CRPS	0.200 ±0.011	0.200 ±0.011	0.199 ±0.011	0.209 ±0.011	0.210 ±0.011	0.260 ±0.011	0.274 ±0.012	0.211 ±0.011	0.210 ±0.011	0.272 ±0.011	0.187 ±0.009	0.187 ±0.009	0.186 ±0.009	0.187 ±0.009	0.337 ±0.005	0.199 ±0.008
	Log	0.195 ±0.011	0.195 ±0.011	0.197 ±0.011	0.203 ±0.010	-	-	0.276 ±0.012	-	0.210 ±0.011	0.272 ±0.011	0.195 ±0.008	0.195 ±0.008	0.195 ±0.008	0.196 ±0.008	-	-
	SE	0.200 ±0.011	0.200 ±0.011	0.200 ±0.011	0.210 ±0.011	0.210 ±0.011	0.272 ±0.011	0.272 ±0.011	0.210 ±0.011	0.210 ±0.011	0.272 ±0.011	0.187 ±0.011	0.187 ±0.011	0.187 ±0.011	0.187 ±0.011	-	-
	Quad.	0.199 ±0.012	0.199 ±0.012	0.196 ±0.011	0.204 ±0.011	0.209 ±0.011	0.247 ±0.011	0.278 ±0.012	0.214 ±0.012	0.210 ±0.011	0.272 ±0.011	0.217 ±0.008	0.217 ±0.008	0.216 ±0.008	0.216 ±0.008	0.462 ±0.005	0.241 ±0.008

from CIFAR-10 (Krizhevsky & Hinton, 2009), SVHN (Netzer et al., 2011), or Fashion-MNIST (Xiao et al., 2017), as well as a mixture of these three sources. We also perform a fine-grained analysis of positional effects by replacing exactly one quadrant at a time with EMNIST (Cohen et al., 2017) (upper-left, upper-right, bottom-left, or bottom-right), keeping the others from MNIST.

**Training details.** Across experiments, we use the same small CNN as in the selective prediction setup, with a larger first linear layer to accommodate the increased input size. We keep hyperparameters

Table 10: Out-of-distribution detection results (AUROC $\uparrow$ ). The ID dataset is MNIST mosaic, where four MNIST digits are tiled in a grid with each corresponding to a digit of the target value. OOD datasets use the corresponding datasets as sources for the tiles. The EMNIST variations replace just one tile with a letter from the EMNIST dataset.

$\mathcal{D}_{\text{OOD}}$	SR	$\hat{r}_{\text{Tot}}^{1,1}$	$\hat{r}_{\text{Tot}}^{2,1}$	$\hat{r}_{\text{Tot}}^{3a,1}$	$\hat{r}_{\text{Tot}}^{3b,1}$	$\hat{r}_{\text{Tot}}^{3a,2}$	$\hat{r}_{\text{Tot}}^{3b,2}$	$\hat{r}_{\text{Bayes}}^1$	$\hat{r}_{\text{Bayes}}^2$	$\hat{r}_{\text{Bayes}}^{3a}$	$\hat{r}_{\text{Bayes}}^{3b}$	$\hat{r}_{\text{Exc}}^{1,1}$	$\hat{r}_{\text{Exc}}^{2,1}$	$\hat{r}_{\text{Exc}}^{3a,1}$	$\hat{r}_{\text{Exc}}^{3b,1}$	$\hat{r}_{\text{Exc}}^{3a,2}$	$\hat{r}_{\text{Exc}}^{3b,2}$	
CIFAR-10	CRPS	0.969 $\pm 0.18$	0.969 $\pm 0.18$	0.969 $\pm 0.18$	0.953 $\pm 0.22$	0.864 $\pm 0.06$	0.831 $\pm 0.13$	0.949 $\pm 0.21$	0.949 $\pm 0.15$	0.832 $\pm 0.21$	0.832 $\pm 0.15$	0.984 $\pm 0.11$	0.984 $\pm 0.11$	0.984 $\pm 0.11$	0.985 $\pm 0.11$	0.963 $\pm 0.16$	0.984 $\pm 0.12$	
	Log	0.977 $\pm 0.16$	0.977 $\pm 0.16$	0.972 $\pm 0.17$	0.962 $\pm 0.22$	-	-	0.830 $\pm 0.15$	0.949 $\pm 0.21$	0.832 $\pm 0.21$	0.832 $\pm 0.15$	0.984 $\pm 0.12$	0.984 $\pm 0.12$	0.984 $\pm 0.13$	0.984 $\pm 0.13$	-	-	
	SE	0.966 $\pm 0.18$	0.966 $\pm 0.18$	0.966 $\pm 0.18$	0.949 $\pm 0.21$	0.949 $\pm 0.21$	0.832 $\pm 0.15$	0.832 $\pm 0.15$	0.949 $\pm 0.21$	0.949 $\pm 0.21$	0.832 $\pm 0.15$	0.832 $\pm 0.15$	0.982 $\pm 0.11$	0.982 $\pm 0.11$	0.982 $\pm 0.11$	0.982 $\pm 0.11$	-	-
	Quad.	0.976 $\pm 0.17$	0.976 $\pm 0.17$	0.976 $\pm 0.17$	0.962 $\pm 0.22$	0.950 $\pm 0.21$	0.893 $\pm 0.18$	0.829 $\pm 0.16$	0.949 $\pm 0.21$	0.949 $\pm 0.21$	0.832 $\pm 0.15$	0.832 $\pm 0.15$	0.972 $\pm 0.22$	0.972 $\pm 0.22$	0.972 $\pm 0.22$	0.974 $\pm 0.21$	0.950 $\pm 0.19$	0.980 $\pm 0.17$
SVHN	CRPS	0.984 $\pm 0.11$	0.984 $\pm 0.11$	0.983 $\pm 0.11$	0.978 $\pm 0.13$	0.938 $\pm 0.12$	0.919 $\pm 0.12$	0.976 $\pm 0.13$	0.976 $\pm 0.12$	0.920 $\pm 0.12$	0.920 $\pm 0.12$	0.986 $\pm 0.12$	0.986 $\pm 0.12$	0.986 $\pm 0.12$	0.986 $\pm 0.12$	0.964 $\pm 0.22$	0.983 $\pm 0.16$	
	Log	0.986 $\pm 0.11$	0.986 $\pm 0.11$	0.985 $\pm 0.11$	0.981 $\pm 0.13$	-	-	0.919 $\pm 0.13$	0.976 $\pm 0.12$	0.920 $\pm 0.12$	0.920 $\pm 0.12$	0.980 $\pm 0.18$	0.980 $\pm 0.18$	0.980 $\pm 0.19$	0.979 $\pm 0.19$	-	-	
	SE	0.982 $\pm 0.11$	0.982 $\pm 0.11$	0.982 $\pm 0.11$	0.976 $\pm 0.12$	0.976 $\pm 0.12$	0.920 $\pm 0.12$	0.920 $\pm 0.12$	0.976 $\pm 0.12$	0.976 $\pm 0.12$	0.920 $\pm 0.12$	0.976 $\pm 0.12$	0.987 $\pm 0.10$	0.987 $\pm 0.10$	0.987 $\pm 0.10$	0.987 $\pm 0.10$	-	-
	Quad.	0.986 $\pm 0.11$	0.986 $\pm 0.11$	0.986 $\pm 0.11$	0.981 $\pm 0.14$	0.976 $\pm 0.14$	0.953 $\pm 0.10$	0.918 $\pm 0.13$	0.976 $\pm 0.13$	0.976 $\pm 0.12$	0.920 $\pm 0.12$	0.920 $\pm 0.12$	0.950 $\pm 0.45$	0.950 $\pm 0.45$	0.950 $\pm 0.45$	0.954 $\pm 0.45$	0.943 $\pm 0.34$	0.968 $\pm 0.31$
Fashion-MNIST	CRPS	0.936 $\pm 0.14$	0.936 $\pm 0.14$	0.936 $\pm 0.14$	0.911 $\pm 0.17$	0.908 $\pm 0.18$	0.825 $\pm 0.25$	0.804 $\pm 0.28$	0.908 $\pm 0.17$	0.805 $\pm 0.18$	0.805 $\pm 0.28$	0.971 $\pm 0.10$	0.971 $\pm 0.10$	0.971 $\pm 0.10$	0.972 $\pm 0.10$	0.936 $\pm 0.05$	0.973 $\pm 0.10$	
	Log	0.950 $\pm 0.12$	0.950 $\pm 0.12$	0.941 $\pm 0.13$	0.922 $\pm 0.16$	-	-	0.804 $\pm 0.28$	0.907 $\pm 0.18$	0.805 $\pm 0.28$	0.805 $\pm 0.28$	0.974 $\pm 0.11$	0.974 $\pm 0.11$	0.974 $\pm 0.11$	0.973 $\pm 0.11$	-	-	
	SE	0.932 $\pm 0.14$	0.932 $\pm 0.14$	0.932 $\pm 0.14$	0.907 $\pm 0.18$	0.907 $\pm 0.18$	0.805 $\pm 0.28$	0.805 $\pm 0.28$	0.907 $\pm 0.18$	0.907 $\pm 0.18$	0.805 $\pm 0.28$	0.966 $\pm 0.10$	0.966 $\pm 0.10$	0.966 $\pm 0.10$	0.966 $\pm 0.10$	-	-	
	Quad.	0.946 $\pm 0.12$	0.946 $\pm 0.12$	0.947 $\pm 0.13$	0.921 $\pm 0.16$	0.908 $\pm 0.17$	0.844 $\pm 0.25$	0.803 $\pm 0.28$	0.907 $\pm 0.17$	0.907 $\pm 0.17$	0.805 $\pm 0.28$	0.963 $\pm 0.17$	0.963 $\pm 0.17$	0.963 $\pm 0.17$	0.964 $\pm 0.17$	0.922 $\pm 0.06$	0.970 $\pm 0.10$	
Mixture (C,S,F)	CRPS	0.963 $\pm 0.13$	0.963 $\pm 0.13$	0.963 $\pm 0.13$	0.947 $\pm 0.15$	0.944 $\pm 0.15$	0.875 $\pm 0.10$	0.851 $\pm 0.17$	0.944 $\pm 0.15$	0.852 $\pm 0.15$	0.852 $\pm 0.17$	0.981 $\pm 0.10$	0.981 $\pm 0.10$	0.981 $\pm 0.10$	0.981 $\pm 0.10$	0.955 $\pm 0.14$	0.980 $\pm 0.12$	
	Log	0.971 $\pm 0.12$	0.971 $\pm 0.12$	0.966 $\pm 0.13$	0.955 $\pm 0.15$	-	-	0.850 $\pm 0.17$	0.944 $\pm 0.15$	0.852 $\pm 0.15$	0.852 $\pm 0.17$	0.979 $\pm 0.13$	0.979 $\pm 0.13$	0.979 $\pm 0.13$	0.979 $\pm 0.13$	-	-	
	SE	0.960 $\pm 0.12$	0.960 $\pm 0.12$	0.960 $\pm 0.12$	0.944 $\pm 0.13$	0.944 $\pm 0.13$	0.852 $\pm 0.10$	0.852 $\pm 0.17$	0.944 $\pm 0.15$	0.944 $\pm 0.15$	0.852 $\pm 0.17$	0.979 $\pm 0.13$	0.979 $\pm 0.13$	0.979 $\pm 0.13$	0.979 $\pm 0.13$	-	-	
	Quad.	0.969 $\pm 0.12$	0.969 $\pm 0.12$	0.970 $\pm 0.12$	0.954 $\pm 0.15$	0.945 $\pm 0.15$	0.896 $\pm 0.13$	0.849 $\pm 0.17$	0.944 $\pm 0.15$	0.944 $\pm 0.15$	0.852 $\pm 0.17$	0.962 $\pm 0.28$	0.962 $\pm 0.28$	0.962 $\pm 0.28$	0.964 $\pm 0.27$	0.938 $\pm 0.19$	0.973 $\pm 0.20$	
EMNIST (tl)	CRPS	0.836 $\pm 0.12$	0.836 $\pm 0.12$	0.836 $\pm 0.12$	0.800 $\pm 0.12$	0.797 $\pm 0.13$	0.718 $\pm 0.13$	0.709 $\pm 0.14$	0.797 $\pm 0.12$	0.796 $\pm 0.12$	0.709 $\pm 0.14$	0.915 $\pm 0.12$	0.915 $\pm 0.12$	0.915 $\pm 0.12$	0.915 $\pm 0.12$	0.846 $\pm 0.20$	0.919 $\pm 0.12$	
	Log	0.855 $\pm 0.12$	0.855 $\pm 0.12$	0.842 $\pm 0.12$	0.811 $\pm 0.12$	-	-	0.709 $\pm 0.14$	-	0.796 $\pm 0.12$	0.709 $\pm 0.14$	0.922 $\pm 0.17$	0.922 $\pm 0.17$	0.922 $\pm 0.17$	0.921 $\pm 0.17$	-	-	
	SE	0.831 $\pm 0.12$	0.831 $\pm 0.12$	0.831 $\pm 0.12$	0.796 $\pm 0.12$	0.796 $\pm 0.12$	0.709 $\pm 0.13$	0.709 $\pm 0.14$	0.796 $\pm 0.12$	0.796 $\pm 0.12$	0.709 $\pm 0.14$	0.903 $\pm 0.12$	0.903 $\pm 0.12$	0.903 $\pm 0.12$	0.903 $\pm 0.12$	-	-	
	Quad.	0.850 $\pm 0.12$	0.850 $\pm 0.12$	0.850 $\pm 0.12$	0.809 $\pm 0.12$	0.797 $\pm 0.12$	0.731 $\pm 0.13$	0.708 $\pm 0.14$	0.797 $\pm 0.12$	0.796 $\pm 0.12$	0.709 $\pm 0.14$	0.910 $\pm 0.11$	0.910 $\pm 0.11$	0.910 $\pm 0.11$	0.911 $\pm 0.11$	0.827 $\pm 0.26$	0.919 $\pm 0.12$	
EMNIST (tr)	CRPS	0.559 $\pm 0.10$	0.559 $\pm 0.10$	0.559 $\pm 0.10$	0.550 $\pm 0.09$	0.550 $\pm 0.09$	0.534 $\pm 0.07$	0.533 $\pm 0.07$	0.550 $\pm 0.09$	0.533 $\pm 0.07$	0.533 $\pm 0.07$	0.594 $\pm 0.14$	0.594 $\pm 0.14$	0.594 $\pm 0.14$	0.594 $\pm 0.14$	0.580 $\pm 0.10$	0.595 $\pm 0.13$	
	Log	0.563 $\pm 0.11$	0.563 $\pm 0.11$	0.561 $\pm 0.11$	0.552 $\pm 0.09$	-	-	0.533 $\pm 0.07$	-	0.550 $\pm 0.09$	0.533 $\pm 0.07$	0.600 $\pm 0.16$	0.600 $\pm 0.16$	0.600 $\pm 0.16$	0.599 $\pm 0.16$	-	-	
	SE	0.558 $\pm 0.10$	0.558 $\pm 0.10$	0.558 $\pm 0.10$	0.550 $\pm 0.09$	0.550 $\pm 0.09$	0.533 $\pm 0.07$	0.533 $\pm 0.07$	0.550 $\pm 0.09$	0.550 $\pm 0.09$	0.533 $\pm 0.07$	0.589 $\pm 0.12$	0.589 $\pm 0.12$	0.589 $\pm 0.12$	0.589 $\pm 0.12$	-	-	
	Quad.	0.562 $\pm 0.10$	0.562 $\pm 0.10$	0.562 $\pm 0.10$	0.552 $\pm 0.09$	0.550 $\pm 0.09$	0.536 $\pm 0.07$	0.533 $\pm 0.07$	0.550 $\pm 0.09$	0.550 $\pm 0.09$	0.533 $\pm 0.07$	0.598 $\pm 0.16$	0.598 $\pm 0.16$	0.599 $\pm 0.16$	0.598 $\pm 0.16$	0.578 $\pm 0.09$	0.597 $\pm 0.14$	
EMNIST (bl)	CRPS	0.557 $\pm 0.12$	0.557 $\pm 0.12$	0.557 $\pm 0.12$	0.550 $\pm 0.11$	0.549 $\pm 0.11$	0.535 $\pm 0.19$	0.534 $\pm 0.18$	0.549 $\pm 0.21$	0.534 $\pm 0.18$	0.534 $\pm 0.18$	0.591 $\pm 0.16$	0.591 $\pm 0.16$	0.591 $\pm 0.16$	0.591 $\pm 0.16$	0.561 $\pm 0.13$	0.593 $\pm 0.21$	
	Log	0.560 $\pm 0.12$	0.560 $\pm 0.12$	0.558 $\pm 0.12$	0.551 $\pm 0.11$	-	-	0.534 $\pm 0.18$	-	0.549 $\pm 0.21$	0.534 $\pm 0.18$	0.595 $\pm 0.16$	0.595 $\pm 0.16$	0.594 $\pm 0.16$	0.594 $\pm 0.16$	-	-	
	SE	0.556 $\pm 0.12$	0.556 $\pm 0.12$	0.556 $\pm 0.12$	0.549 $\pm 0.11$	0.549 $\pm 0.11$	0.534 $\pm 0.18$	0.534 $\pm 0.18$	0.549 $\pm 0.21$	0.549 $\pm 0.21$	0.534 $\pm 0.18$	0.586 $\pm 0.22$	0.586 $\pm 0.22$	0.586 $\pm 0.22$	0.586 $\pm 0.22$	-	-	
	Quad.	0.560 $\pm 0.12$	0.560 $\pm 0.12$	0.560 $\pm 0.12$	0.551 $\pm 0.11$	0.549 $\pm 0.11$	0.537 $\pm 0.19$	0.534 $\pm 0.18$	0.549 $\pm 0.21$	0.549 $\pm 0.21$	0.534 $\pm 0.18$	0.593 $\pm 0.18$	0.593 $\pm 0.18$	0.593 $\pm 0.18$	0.593 $\pm 0.18$	0.557 $\pm 0.14$	0.595 $\pm 0.19$	
EMNIST (br)	CRPS	0.548 $\pm 0.13$	0.548 $\pm 0.13$	0.548 $\pm 0.13$	0.541 $\pm 0.14$	0.541 $\pm 0.13$	0.529 $\pm 0.13$	0.528 $\pm 0.13$	0.541 $\pm 0.13$	0.528 $\pm 0.13$	0.528 $\pm 0.13$	0.577 $\pm 0.08$	0.577 $\pm 0.08$	0.577 $\pm 0.08$	0.577 $\pm 0.08$	0.553 $\pm 0.19$	0.577 $\pm 0.09$	
	Log	0.551 $\pm 0.13$	0.551 $\pm 0.13$	0.549 $\pm 0.13$	0.542 $\pm 0.14$	-	-	0.528 $\pm 0.13$	-	0.541 $\pm 0.13$	0.528 $\pm 0.13$	0.581 $\pm 0.09$	0.581 $\pm 0.09$	0.581 $\pm 0.09$	0.581 $\pm 0.09$	-	-	
	SE	0.547 $\pm 0.13$	0.547 $\pm 0.13$	0.547 $\pm 0.13$	0.541 $\pm 0.13$	0.541 $\pm 0.13$	0.528 $\pm 0.13$	0.528 $\pm 0.13$	0.541 $\pm 0.13$	0.541 $\pm 0.13$	0.528 $\pm 0.13$	0.573 $\pm 0.11$	0.573 $\pm 0.11$	0.573 $\pm 0.11$	0.573 $\pm 0.11$	-	-	
	Quad.	0.550 $\pm 0.13$	0.550 $\pm 0.13$	0.550 $\pm 0.13$	0.542 $\pm 0.14$	0.541 $\pm 0.13$	0.530 $\pm 0.13$	0.528 $\pm 0.13$	0.541 $\pm 0.13$	0.541 $\pm 0.13$	0.528 $\pm 0.13$	0.579 $\pm 0.08$	0.579 $\pm 0.08$	0.579 $\pm 0.08$	0.579 $\pm 0.08$	0.552 $\pm 0.20$	0.578 $\pm 0.08$	

identical to the selective prediction setting, and train ensembles of 10 models on the in-distribution MNIST-mosaic task.

**Extended results.** For evaluation, we generate matched in-distribution and OOD sets of equal size and compute AUROC using uncertainty scores to discriminate OOD from in-distribution samples. Full per-dataset results are listed in Table 10. Replacing all four quadrants with another dataset yields similar detection performance across sources (first four rows), with SVHN—despite also containing digits—detected comparably well. In the positional analysis (last four rows), replacing the upper-left digit (the most significant position in the four-digit number

1998 performance (SGEMM; Ballester-Ripoll et al., 2019), Combined Cycle Power Plant (CCPP; Tüfekci,  
1999 2014), Physicochemical Properties of Protein Tertiary Structure (CASP; Rana, 2013), Online News  
2000 Popularity (NEWS; Fernandes et al., 2015) and BlogFeedback (BLOG; Buza, 2013). For all datasets,  
2001 we split 20% into a test set for evaluation. The remaining set was further split into validation and pool  
2002 set. We selected around 200 samples (exact sizes per dataset in code, which was selected based initial  
2003 performance compared to performance on full pool set) as a training set to start the active learning  
2004 experiments.

2005 **Training details.** In each acquisition iteration, we selected a number of additional samples from  
2006 the pool to augment the training set. We do not strictly select based on the score induced by a given  
2007 uncertainty measure, but obtain a categorical distribution based on applying the softmax over the  
2008 uncertainty score for all samples in the pool set. We then sample without replacement according  
2009 to this distribution. This avoids selecting very similar samples within an acquisition batch, which  
2010 is a known issue in practice (Kirsch et al., 2019). For these experiments, we utilized ensembles of  
2011 standard three-layer MLPs with ReLU activations with hidden sizes of 100.

2012 **Results.** The full results are provided in Figure 10. We observe that for most settings, using the  
2013 uncertainty scores as acquisition function improves upon random sampling. Notable exceptions are  
2014 total and Bayes risks for the quadratic score, which perform rather poorly.  
2015

2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

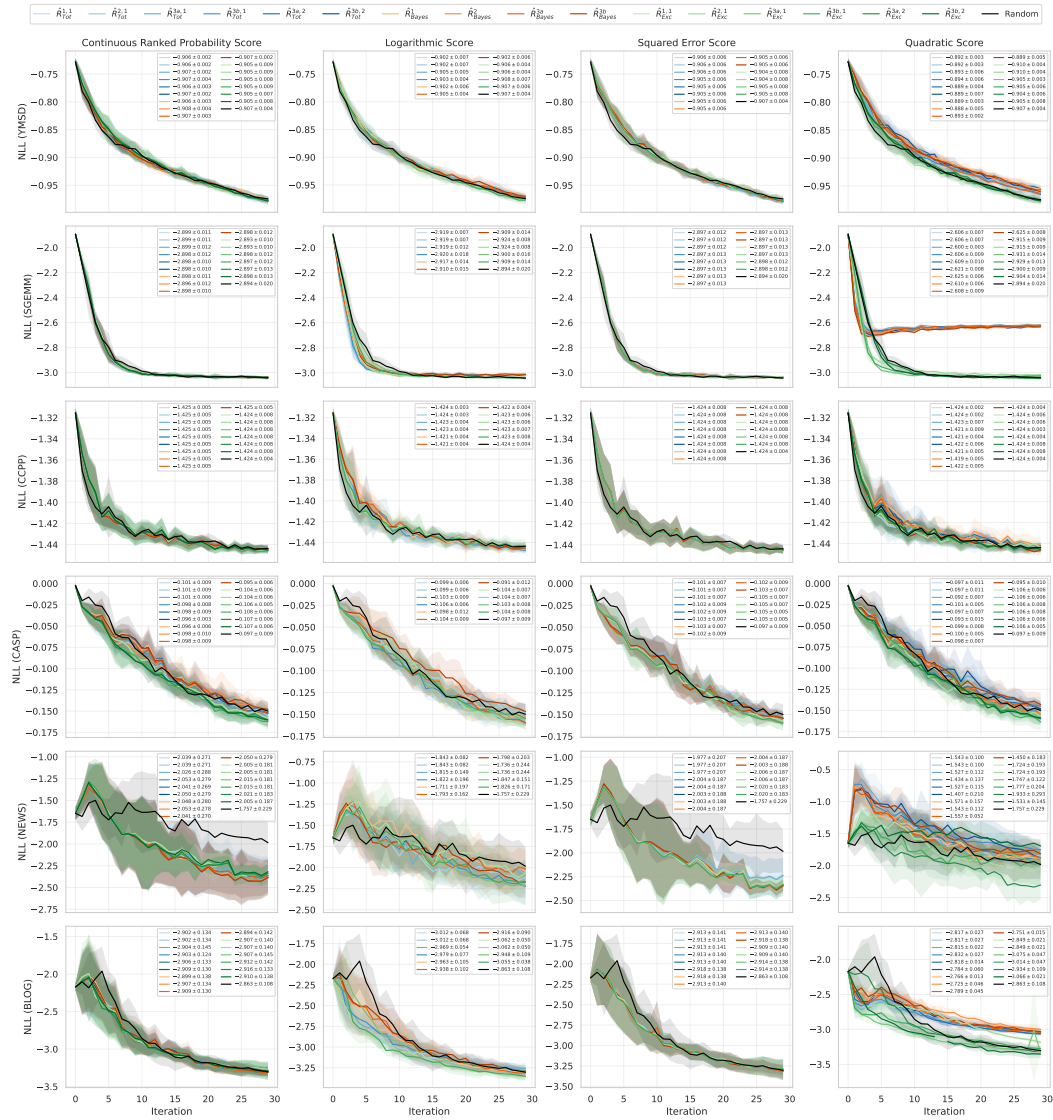


Figure 10: Detailed results on active learning. We report the negative log-likelihood (NLL) of the ensemble after training for each acquisition iteration. In the insets are the average NLLs over the iterations as summary statistic to compare different acquisition functions. Statistics are obtained over five runs.