

# Benchmarking Bayesian neural networks and evaluation metrics for regression tasks

Brian Staber <sup>\*1</sup> Sébastien Da Veiga <sup>2</sup>

## Abstract

Due to the growing adoption of deep neural networks in many fields of science and engineering, modeling and estimating their uncertainties has become of primary importance. Despite the growing literature about uncertainty quantification in deep learning, the quality of the uncertainty estimates remains an open question. In this work, we assess for the first time the performance of several approximation methods for Bayesian neural networks on regression tasks by evaluating the quality of the confidence regions with several coverage metrics. The selected algorithms are also compared in terms of predictivity, kernelized Stein discrepancy and maximum mean discrepancy with respect to a reference posterior in both weight and function space. Our findings show that (i) some algorithms have excellent predictive performance but tend to largely over or underestimate uncertainties (ii) it is possible to achieve good accuracy and a given target coverage with finely tuned hyperparameters and (iii) the promising **kernel Stein discrepancy** cannot be exclusively relied on to assess the posterior approximation. As a by-product of this benchmark, we also compute and visualize the similarity of all algorithms and corresponding hyperparameters: interestingly we identify a few clusters of algorithms with similar behavior in weight space, giving new insights on how they explore the posterior distribution.

## 1. Introduction

Due to their recent achievements in the last decade, deep neural networks have been widely adopted in many research fields and industries. However, they still suffer from several

<sup>\*</sup>Equal contribution <sup>1</sup>Safran Tech, Digital Sciences & Technologies, 78114 Magny-Les-Hameaux, France <sup>2</sup>ENSAI, CREST, F-35000 Rennes, France. Correspondence to: Brian Staber <brian.staber@safrangroup.com>, Sébastien Da Veiga <sebastien.da-veiga@ensai.fr>.

shortcomings that prevent their deployment in fields where decisions involve high stakes. These limitations are mainly due to their inability to provide uncertainty estimates which are crucial for many real-world applications. Uncertainty quantification in deep learning has attracted a lot of attention and several approaches have been investigated. Thorough overviews are provided by the recent review papers of Abdar et al. (2021) and Gawlikowski et al. (2021). Amongst the available approaches, Bayesian neural networks offer a simple and flexible formulation but raise several challenges. The high-dimensional posterior distribution of the network parameters is intractable, possibly multimodal, and the influence of the prior distribution remains an open question. The posterior distribution is typically approximated using sampling methods, variational inference, or Gaussian approximations. When dealing with complex posterior distributions, gradient-based sampling methods such as Markov Chain Monte Carlo (MCMC) methods are often adopted. In particular, Hamiltonian Monte Carlo (HMC) is usually considered as a gold standard sampling algorithm but is unfortunately extremely computationally demanding (Izmailov et al., 2021; Cobb & Jalaian, 2021). Since the work of Welling & Teh (2011), stochastic gradient MCMC methods have been extensively studied (see, e.g., (Ma et al., 2015) and (Nemeth & Fearnhead, 2021)), where the Metropolis-Hastings correction is omitted and a mini-batched stochastic gradient is used. Variational approaches approximate the posterior distribution by minimizing the Kullback-Leibler divergence over a family of tractable distributions (Hoffman et al., 2013). Several scalable methods have been proposed such as the Bayes by Backprop (BBB) method of Graves (2011) and Blundell et al. (2015), the multiplicative normalizing flows proposed by Louizos & Welling (2017), and the Monte Carlo dropout method (Gal & Ghahramani, 2016). Although ensemble methods are generally used to improve the generalization capabilities of neural networks, they can also be interpreted as a Bayesian approach (Lakshminarayanan et al., 2017; Fort et al., 2019), where uncertainties are predicted thanks to random initialization and data shuffling. Gaussian approximations have also been investigated by several authors, such as the Laplace approximation (Ritter et al., 2018; Daxberger et al., 2021) and the highly scalable stochastic weight averaging Gaussian method (Maddox et al., 2019).

Regardless of the approximation method, evaluating the quality of the predictive uncertainties remains an open issue. In contrast to previous works, we propose for the first time to estimate sensible coverage probabilities by taking into account the variability induced by the training dataset. We also compare the selected algorithms with the help of discrepancy measures, namely, the maximum mean discrepancy (Gretton et al., 2006) and the kernelized Stein discrepancy (Liu et al., 2016).

## 2. Related works

Izmailov et al. (2021) have studied the performance of Bayesian neural networks (BNNs) by relying on exhaustive full-batch Hamiltonian Monte Carlo. The performance is assessed in terms of RMSE for regression problems and accuracy for classification tasks. Most notably, Izmailov et al. (2021) have applied BNNs to practical deep architectures and large datasets thanks to impressive computational resources. Wenzel et al. (2020) have studied the influence of posterior temperature on the performance of stochastic gradient sampling methods. The performance is reported in terms of accuracy for classification tasks only. More closely related to this work, Yao et al. (2019) has compared several approximation methods, for regression and classification tasks, and reported the prediction interval coverage probability in order to evaluate the quality of the approximated prediction intervals.

The objective of this work is to assess the performance of Bayesian neural networks on simple regression tasks by comparing several metrics. In contrast to previous works, we assess the validity of the confidence intervals with marginal and conditional coverage probabilities instead of predictive interval coverage probability, which are defined in section 4. It is worth emphasizing that the aim of this work is not to suggest to use marginal and conditional coverage for practical problems, but to investigate which performance BNNs can achieve. In particular, we investigate if there are correlations between accuracy, coverage metrics and kernel Stein discrepancy, as well as if some algorithms have the same exploration behavior in weight and function spaces.

## 3. Bayesian neural networks

Let  $\mathcal{D} = \{(\mathbf{X}_i, \mathbf{Y}_i)\}_{i=1}^N$  denote a dataset made of  $N$  observations of input  $\mathbf{X}_i \in \mathbb{R}^D$  and output  $\mathbf{Y}_i \in \mathbb{R}^M$  pairs. The observations take the form  $\mathbf{Y}_i = f(\mathbf{X}_i) + \varepsilon_i$  where  $f : \mathbb{R}^D \rightarrow \mathbb{R}^M$  represents the unknown latent regression function, and  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2(\mathbf{X}_i)\mathbf{I}_M)$  is an additive noise modeling aleatoric uncertainties. The latent function  $f$  is approximated by a neural network  $\hat{f}(\cdot; \mathbf{w})$  with parameters  $\mathbf{w} \in \mathbb{R}^d$ . The observations  $\mathbf{Y}_1, \dots, \mathbf{Y}_N$  are

assumed to be i.i.d. and the likelihood function takes the form  $p(\mathbf{Y}_i | \mathbf{X}_i, \mathbf{w}) = \mathcal{N}(\mathbf{Y}_i; \hat{f}(\mathbf{X}_i; \mathbf{w}), \sigma^2 \mathbf{I}_M)$  for any  $i \in \{1, \dots, N\}$ . Given a prior distribution  $p(\mathbf{w})$  over the network parameters, the posterior distribution  $p(\mathbf{w} | \mathbf{X}, \mathbf{Y})$  can be deduced using Baye's formula:  $p(\mathbf{w} | \mathbf{X}, \mathbf{Y}) \propto p(\mathbf{Y} | \mathbf{X}, \mathbf{w})p(\mathbf{w})$ . The prediction of the model for a new test input  $\mathbf{x}$  is then given by

$$p(\mathbf{y} | \mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int_{\mathbb{R}^d} p(\mathbf{y} | \mathbf{x}, \mathbf{w})p(\mathbf{w} | \mathbf{X}, \mathbf{Y}) d\mathbf{w}.$$

Unfortunately, this integral is intractable and we must have recourse to approximate inference. The true posterior distribution is approximated by some  $q(\mathbf{w}) \approx p(\mathbf{w} | \mathbf{X}, \mathbf{Y})$  obtained via for instance MCMC methods, variational inference, or Gaussian approximations. In the rest of this section, we briefly describe the approximation methods that we consider in this paper, additional details are given in Appendix A.

**Monte Carlo Markov Chain.** Monte Carlo Markov chain (MCMC) methods generate a Markov chain whose stationary distribution is the target posterior distribution. In this work, we consider the **Hamiltonian Monte Carlo (HMC)** algorithm (Neal, 2011) which requires the knowledge of the unnormalized posterior  $p(\mathbf{Y} | \mathbf{X}, \mathbf{w})p(\mathbf{w})$  and the gradient of the potential function  $U(\mathbf{w}) = -\log(p(\mathbf{Y} | \mathbf{X}, \mathbf{w})p(\mathbf{w}))$ , given by

$$\nabla U(\mathbf{w}) = - \sum_{i=1}^N \nabla U_i(\mathbf{w}) - \nabla \log(p(\mathbf{w})),$$

where  $U_i(\mathbf{w}) = \log(p(\mathbf{Y}_i | \mathbf{X}_i, \mathbf{w}))$ .

**Stochastic gradient Monte Carlo Markov Chain.** In order to alleviate the computational cost of classical MCMC methods, many efforts have been dedicated to the development of stochastic gradient MCMC methods (Welling & Teh, 2011; Ahn et al., 2012; Ding et al., 2014; Chen et al., 2014; Ma et al., 2015; Li et al., 2016; Zhang et al., 2019). Here, the Metropolis-Hastings correction step is omitted and the gradient of the potential function is approximated by a stochastic, mini-batched, gradient:

$$\hat{\nabla} U(\mathbf{w}) = - \frac{N}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla U_i(\mathbf{w}) - \nabla \log(p(\mathbf{w})).$$

Despite being computationally more efficient than traditional MCMC, stochastic gradient MCMC methods introduce asymptotic bias. In this work, we consider the stochastic gradient Lagenvin dynamics (**SGLD**) and Hamiltonian Monte Carlo (**SGHMC**). We also consider variants of SGMCMC methods that include variance reduction techniques (Baker et al., 2019; Dubey et al., 2016), where the stochastic gradient  $\hat{\nabla} U$  is replaced by the following estimation:

$$\tilde{\nabla} U(\mathbf{w}) = \nabla U(\mathbf{w})|_{\mathbf{w}=\eta} + \hat{\nabla} U(\mathbf{w}) - \hat{\nabla} U(\mathbf{w})|_{\mathbf{w}=\eta}.$$

We then consider the variants **SGLD-CV** and **SGHMC-CV** where  $\eta$  is set to a MAP estimate  $\mathbf{w}_{\text{MAP}}$ , together with the variants **SGLD-SVRG** and **SGHMC-SVRG** where  $\eta$  is first set to a MAP estimate, and then updated to the current value of  $\mathbf{w}$  every  $m$  iterations. Finally, we also evaluate the preconditioned method **pSGLD** of Li et al. (2016), and the cyclical algorithms **C-SGLD** and **C-SGHMC** that rely on a cyclical step size (Zhang et al., 2019).

**Gaussian approximations.** We consider two Gaussian approximations: **SWAG** and the Laplace approximation with a kronecker factored log likelihood Hessian approximation (**LA-KFAC**), which both require a pre-trained neural network with parameters  $\mathbf{w}_{\text{MAP}}$ . The **SWAG** approximation is constructed by collecting values of the parameters along a SGD trajectory with a possibly high step size. The **SWAG** approximation reads as  $p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) \approx \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{SWAG}}, \Sigma_{\text{SWAG}})$  where  $\mathbf{w}_{\text{SWAG}}$  is given by the running mean of the SGD iterates, and  $\Sigma_{\text{SWAG}}$  is a diagonal plus low rank approximation. The **LA-KFAC** approximation is given by  $p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) \approx \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \Sigma_{\text{KFAC}})$  where  $\Sigma_{\text{KFAC}}$  denotes the Kronecker-factored approximate curvature approximation (see, e.g., Daxberger et al. (2021)).

**Variational methods.** The **Monte Carlo dropout method (MC-Dropout)** proposed by Gal & Ghahramani (2016) is considered herein. The neural network is augmented with dropout layers which are activated during both training and inference stages. The dropout mechanism is here applied to the output features of each layer.

**Deep ensembles.** We also investigate the **deep ensembles** method of Lakshminarayanan et al. (2017); Fort et al. (2019), which consists in training several neural networks independently with random initializations, and gathering their predictions to obtain a mean prediction and uncertainties.

## 4. Evaluation metrics

**Validity of the confidence intervals.** We assess the validity of the confidence intervals produced by an approximation method with coverage probabilities. There are several related notions of coverage probabilities such as the prediction interval coverage probability, marginal coverage probability, and the conditional coverage probability (Lin et al., 2021). Given a target confidence level  $1 - \alpha$ , the confidence interval  $\hat{\mathcal{C}}_\alpha^{\mathcal{D}}$  is said to have a valid *prediction interval coverage probability* (PICP) if

$$\mathbb{P}\{\mathbf{Y}^* \in \hat{\mathcal{C}}_\alpha^{\mathcal{D}}(\mathbf{X}^*) | \mathcal{D}\} \geq 1 - \alpha,$$

where the probability is taken only over a test data  $(\mathbf{X}^*, \mathbf{Y}^*)$ . This coverage probability does not take into account the variability induced by the training dataset  $\mathcal{D}$ . A more

suitably property would be the *marginal coverage probability* (MCP):

$$\mathbb{P}\{\mathbf{Y}^* \in \hat{\mathcal{C}}_\alpha^{\mathcal{D}}(\mathbf{X}^*)\} \geq 1 - \alpha,$$

where the probability is taken over to both the training data  $\mathcal{D}$  and the test data  $(\mathbf{X}^*, \mathbf{Y}^*)$ . An even stronger property can be obtained by conditioning on the test input data  $\mathbf{X}^*$ , leading to the *conditional coverage probability* (CCP):

$$\mathbb{P}\{\mathbf{Y}^* \in \hat{\mathcal{C}}_\alpha^{\mathcal{D}}(\mathbf{X}^*) | \mathbf{X}^* = \mathbf{x}\} \geq 1 - \alpha,$$

for almost all  $\mathbf{x} \in \mathcal{X}$ , where the probability is taken over the training dataset  $\mathcal{D}$ .

**Distance to the HMC reference (weight and function space).** We compute the maximum mean discrepancy (MMD) between each approximation and the approximation obtained via exhaustive **HMC**. Let  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be a kernel function and let  $\mathcal{H}(k)$  be the reproducing kernel Hilbert space with kernel  $k$ . The MMD between two probability measures  $\mathbb{Q}$  and  $\mathbb{Q}'$  is defined as (Gretton et al., 2006):

$$\text{MMD}(\mathbb{Q}, \mathbb{Q}') = \|\mu_{\mathbb{Q}} - \mu_{\mathbb{Q}'}\|_{\mathcal{H}(k)}, \quad \mu_{\mathbb{Q}} = \int k(\cdot, \mathbf{w}) d\mathbb{Q},$$

where  $\mu_{\mathbb{Q}}$  is called the kernel mean embedding of  $\mathbb{Q}$  in  $\mathcal{H}(k)$ . Here  $k$  is chosen as the distance-based kernel function  $k(\mathbf{w}, \mathbf{w}') = \|\mathbf{w}\|_2 + \|\mathbf{w}'\|_2 - \|\mathbf{w} - \mathbf{w}'\|_2$  proposed by Sejdinovic et al. (2013). In this setting,  $\text{MMD}(\mathbb{Q}, \mathbb{Q}') = 0$  implies that  $\mathbb{Q} = \mathbb{Q}'$ . For a given training dataset  $\mathcal{D}$  and a set of hyperparameters, we compute the MMD between each approximation  $\hat{\mathbb{Q}}$  in weight space obtained by **SGCM-CMC**, **SWAG**, **LA-KFAC**, **MC-Dropout**, or **Deep ensembles**, and the approximation  $\mathbb{Q}_{\text{HMC}}$  obtained via exhaustive **HMC** that is considered as the reference. In addition, given a sample  $\{\mathbf{w}_i\}_{i=1}^m$  from an approximation  $\hat{\mathbb{Q}}$ , we can deduce the predictions of the neural network  $\mathbf{x} \mapsto f(\mathbf{x}; \mathbf{w})$  for each algorithm and compute the MMD distance to the HMC reference in function space.

**Distance to the target posterior (weight space).** We also assess the performance of an approximation method by measuring a distance to the target posterior distribution. We rely on the kernelized Stein discrepancy which, in contrast to other distances, only requires the knowledge of the gradient of the log-posterior distribution. The squared KSD between the target posterior measure  $\mathbb{P}$  and any other probability measure  $\mathbb{Q}$  is defined as

$$\text{KSD}^2(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{\mathbf{w} \sim \mathbb{Q}} \mathbb{E}_{\mathbf{w}' \sim \mathbb{Q}} k_p(\mathbf{w}, \mathbf{w}'),$$

where  $k_p$  denotes the Stein kernel given by

$$k_p(\mathbf{w}, \mathbf{w}') = \langle \nabla_{\mathbf{w}}, \nabla_{\mathbf{w}'} k(\mathbf{w}, \mathbf{w}') \rangle + \langle s_p(\mathbf{w}), \nabla_{\mathbf{w}'} k(\mathbf{w}, \mathbf{w}') \rangle + \langle s_p(\mathbf{w}'), \nabla_{\mathbf{w}} k(\mathbf{w}, \mathbf{w}') \rangle + \langle s_p(\mathbf{w}), s_p(\mathbf{w}') \rangle k(\mathbf{w}, \mathbf{w}').$$

Here,  $s_p$  denotes the score function of the target posterior distribution, namely,  $s_p(\mathbf{w}) = \nabla_{\mathbf{w}} \log(p(\mathbf{Y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}))$ . The Stein kernel also depends on an additional kernel function,  $k$ , which has to be carefully chosen. Based on the theoretical results of Gorham & Mackey (2017), the kernel  $k$  is chosen as the inverse multi-quadratic (IMQ) kernel, which is defined as  $k(\mathbf{w}, \mathbf{w}') = (1 + \|\mathbf{w} - \mathbf{w}'\|_{\Gamma})^{-1/2}$  with  $\Gamma = \ell^2 I$ . In this setting, the KSD defines a distance between two probability measures such that  $\text{KSD}(\mathbb{P}, \mathbb{Q}) = 0$  implies that  $\mathbb{P} = \mathbb{Q}$ . The lengthscale  $\ell$  of the IMQ kernel is chosen as the median of the pairwise distances, estimated with a subsample of the entire collection of samples generated by all the approximation methods.

**Similarities between the algorithms (weight and function space).** We also use the MMD in order to establish possible similarities between the algorithms. More precisely, the MMD is computed between each pair  $(\hat{\mathbb{Q}}, \hat{\mathbb{Q}}')$  of approximations obtained with the considered approximation methods and set of hyperparameters. The resulting matrix of pairwise MMD distances is visualized thanks to multidimensional scaling (Torgerson, 1952) and analyzed to highlight any similarities between the approximation methods. We also proceed the same in function space.

## 5. Experimental setup

The performances of the selected algorithms are studied for several regression tasks. In this section, we summarize the considered datasets and neural network architectures. For each algorithm, we also study the influence of some hyperparameters, which are summarized in the sequel of this section. Additional details about the experiment setup have been reported in the Supplementary Material.

**Regression tasks.** We consider 4 synthetic regression problems where the output is one-dimensional but the input may be one or multi-dimensional. For each synthetic regression problem, we generate  $N_{\mathcal{D}} = 500$  independent datasets  $\mathcal{D}_1, \dots, \mathcal{D}_{N_{\mathcal{D}}}$  used for MCP and CCP computations. We also consider an additional test dataset  $\mathcal{D}^* = \{(\mathbf{X}_i^*, Y_i^*)\}_{i=1}^{N^*}$  for accuracy, PICP, MCP and CCP. For some problems it may have out-of-distribution (OOD) samples. Each element  $(\mathbf{X}_i, Y_i)$  of a training dataset  $\mathcal{D}_j$  is such that  $Y_i = f(\mathbf{X}_i) + \epsilon_i$ , where  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ ,  $i = 1, \dots, N$ . The underlying latent regression function  $f$  and the variance of the noise  $\sigma$  are both known. The four regression problems are constructed with analytical functions, which are summarized in Table 1 below. More details about the regression problems and their generative process can be found in Appendix B.

**Neural network architectures.** A feed-forward neural network with ReLU activations is used as a surrogate model

Table 1. Description of the synthetic datasets where AF stands for analytical functions.

Task	Latent function	$\sigma$	$D$	$N$	$N^*$	OOD
AF#1	$\cos(2x) + \sin(x)$	0.2	1	100	200	✗
AF#2	$0.1x^2$	0.25	1	100	200	✓
AF#3	$-(1+x)\sin(1.2x)$	0.25	1	82	200	✓
AF#4	$\text{MLP}(\cdot; \mathbf{w}), \mathbf{w} \sim \mathcal{N}(0, \mathbf{I})$	0.02	2	120	120	✓

for every regression problem. The number of parameters ranges from 2,651 to 20,501 (see Appendix C for more details). We use a centered normalized Gaussian prior distribution for the weights in all the experiments. Note that for each training data  $\mathcal{D}_j$ ,  $j = 1, \dots, N_{\mathcal{D}}$ , a MAP estimate is computed by training the neural network with an Adam optimizer and an exponentially decaying learning rate. These  $N_{\mathcal{D}}$  MAP estimates are subsequently used in **LA-KFAC**, **SWAG**, **SGMCMC-CV**, and **SGMCMC-SVRG**.

**Hyperparameters.** Besides the **LA-KFAC** approximation, all the algorithms depend on one or several hyperparameters. For all the remaining methods (**SGMCMC**, **Deep ensembles**, **MC-Dropout**, and **SWAG**), we study the influence of the step size  $\epsilon > 0$  (*i.e.*, the learning rate) by considering 10 equally log-spaced values  $\epsilon_1 < \dots < \epsilon_{10}$ . The ranges are carefully chosen for each type of algorithm and are reported in Appendix A. We also investigate the influence of two additional hyperparameters. Five dropout rates in the **MC-Dropout** method are considered (0.1, 0.2, ..., 0.5) and three cycle lengths in the **C-SGLD** and **C-SGHMC** methods are investigated (10, 100, and 1000). Other hyperparameters are held fixed, see Appendix A.3 for more details.

**HMC reference.** For a given training dataset  $\mathcal{D}$ , a reference sample is generated by Hamiltonian Monte Carlo and subsequently used to evaluate the performance of the selected algorithms (see section 4). We run 3 HMC chains of 200 iterations, discard the first 100 iterations as burn-in, and perform 10,000 leapfrogs steps. The step size is selected such that the Metropolis-Hastings accept rates are at least above 80%.

## 6. Regression without OOD testing

We first present the results for the first experiment (AF#1 in Table 1). Here, the training and test datasets are sampled from the same distribution so that any test input lies within the range of the input training data points.

**Coverage probabilities.** Figure 1 gathers the graphs of the marginal coverage probabilities and the mean absolute error in conditional coverage with respect to the step size  $\epsilon$  for a target coverage of 0.95. In addition, the coefficient

of determination on the testing set, denoted by  $Q^2 = 1 - \sum_{i=1}^n (Y_i^* - \hat{f}(\mathbf{X}_i^*; \mathbf{w}))^2 / \text{Var } Y^*$ , is also reported. It can be observed that

- (a) The **SGLD** and **SGHMC** methods yield similar performances.
- (b) The performance of the variants with control variates (**SGLD-CV** and **SGHMC-CV**) drops significantly for high step sizes. As pointed out by (Nemeth & Fearnhead, 2021), if the state of the Markov chain gets far from the control variate (here a MAP estimate), then the variance of the stochastic gradient can increase instead of being reduced, hence explaining the observed behavior.
- (c) Updating the control variates seems to address this issue as illustrated by the performances of **SGLD-SVRG**, but **SGHMC-SVRG** is still affected by high step sizes. Given that **SGHMC** mixes better than **SGLD**, we believe that this behavior can be corrected by increasing the frequency at which the control variates are updated in **SGHMC-SVRG**.
- (d) The preconditioned SGLD method easily reaches the target coverage probability but has lower coefficients of determination than the other methods.
- (e) **Deep ensembles** and **SWAG** both yield high coefficient of determinations but only **deep ensembles** is able to reach the target coverage probability.
- (f) **MC-dropout** has the lowest performances in this example. In particular, the high variability of the absolute error in conditional coverage suggests that the confidence intervals are not smooth with respect to the input space.
- (g) Finally, for the **LA-KFAC** method we obtain  $Q^2 = 0.99 \pm 0.004$ , MCP =  $0.99 \pm 0.01$ , and MAE =  $0.048 \pm 0.003$ . As a result, the mean prediction of **LA-KFAC** is very accurate but the high coverage probability suggests that it overestimates the uncertainties.

We also investigate the marginal coverage probability with respect to various target levels  $(1 - \alpha)$ ,  $\alpha \in ]0, 1[$ , which are shown in Figure 2 for 6 selected algorithms (see also Appendix C). Here, we see that **pSGLD** and **LA-KFAC** easily overestimate the target level regardless of the step size. Finally, we report the best marginal coverage probabilities (closest to 0.95) and the best  $Q^2$  coefficients (closest to 1) for each algorithm in Table 2. For each best coverage (resp.  $Q^2$ ), we also report the associated  $Q^2$  (resp. coverage). Interestingly, we observe that the best coverages lying in  $0.95 \pm 0.1$  do not especially correspond to the best regression

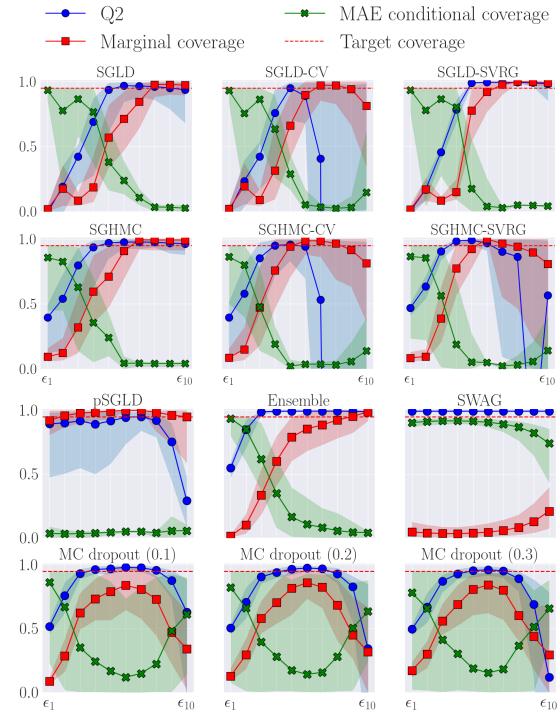


Figure 1. Problem AF#1. Coverage metrics and  $Q^2$  coefficient with respect to the step size  $\epsilon$ . The target coverage is set to 0.95. Results obtained for the cyclical **SGMCMC** variants are reported in the Appendix.

coefficients  $Q^2$ . Most notably, **Deep ensembles** is able to achieve the best performances in order of both coverage and prediction accuracy.

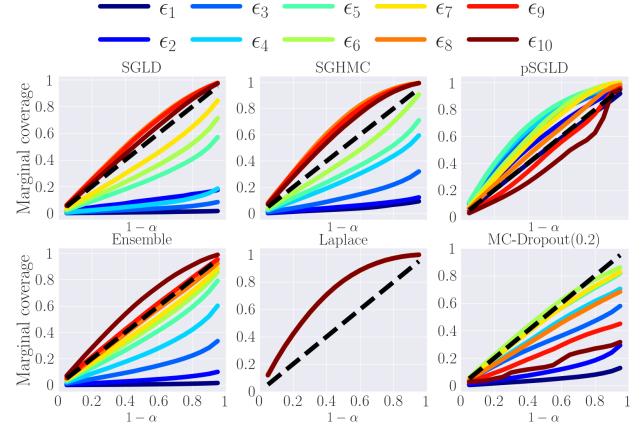


Figure 2. Problem AF#1. Graphs of the marginal coverage probability with respect to the target level and for 10 step sizes.

**MMD distance to the HMC reference.** The distances to the HMC reference sample are shown in Figure 3. In weight space, the MMD distance increases with the step size for

Table 2. Best marginal coverage probabilities (MCP) and best  $Q^2$  coefficients obtained with each algorithm.

Method	Best MCP	$Q^2$	Best $Q^2$ MCP	Method	Best MCP	$Q^2$	Best $Q^2$ MCP		
SGLD	0.97	0.93	0.96	0.71	CSGHMC(10)	0.97	0.97	0.97	0.89
SGLD-CV	<b>0.94</b>	$-10^3$	0.95	0.66	CSGHMC(100)	<b>0.94</b>	0.97	0.97	0.94
SGLD-SVRG	0.97	0.99	<b>0.99</b>	0.91	CSGHMC(1000)	<b>0.94</b>	0.96	0.97	0.93
SGHMC	0.99	0.96	0.97	0.99	pSGLD	<b>0.94</b>	0.29	0.95	0.99
SGHMC-CV	<b>0.94</b>	0.95	0.95	0.94	Deep ensemble	<b>0.95</b>	0.99	<b>0.99</b>	0.88
SGHMC-SVRG	<b>0.94</b>	0.86	<b>0.99</b>	0.92	SWAG	0.20	0.99	<b>0.99</b>	0.03
CSGLD(10)	<b>0.95</b>	0.96	0.96	0.66	MC Drop.(0.1)	0.83	0.98	<b>0.98</b>	0.83
CSGLD(100)	<b>0.94</b>	0.96	0.96	0.76	MC Drop.(0.2)	0.85	0.97	0.97	0.85
CSGLD(1000)	0.92	0.95	0.96	0.75	MC Drop.(0.3)	0.84	0.96	0.96	0.84

most algorithms. **Deep ensembles** is the only method for which the MMD distance in weight space decreases below 5. In contrast, the MMD distances in function space may decrease towards 0. In particular, the **SGLD-SVRG**, **deep ensembles**, and **SWAG** methods yield the lowest MMD distances in function space. With the **LA-KFAC** method, we obtained MMD distances of  $5.76 \pm 0.13$  in weight space, and of  $1.27 \pm 0.21$  in function space. Overall and perhaps surprisingly, the behavior with respect to the step size is not the same in weight and function space.

**KSD distance to the target posterior.** The kernelized Stein discrepancies between the target posterior  $\mathbb{P}$  and the approximations  $\mathbb{Q}$  are shown in Figure 4. Here, the lowest discrepancies are obtained with **deep ensembles** and **SWAG**, while **SGMCMC** methods yield higher KSDs. Surprisingly, the **LA-KFAC** method yields the highest KSDs, with values between  $10^5$  and  $10^7$ . If for some algorithms the conclusions from the MMD and the KSD coincide, this is not the case in general: we discuss below some potential explanations related to KSD limitations.

**Similarities between the algorithms.** The similarities between the algorithms is depicted in Figure 5. At least four groups of approximations can be distinguished in weight space (left panel in Figure 5): **SGLD**, **SGLD-CV** and **CSGLD** methods, **SGHMC**, **SGHMC-CV** and **CSGHMC**, **pSGLD** and **LA-KFAC**, and the remaining methods form the last group. Amongst the **SGLD** and **SGHMC** methods, it is clearly observed that **CV** variants yield distinct approximations. Surprisingly, **SGLD-SVRG** and **SGHMC-SVRG** generate similar approximations to **deep ensembles**, in weight space. The similarities in function space do not exhibit any particular structure and have been reported in Appendix C.2 (see Figure C.2).

**Discussion.** Based on the results presented in this section, a few observations can be made. First, we see that low KSD values do not especially correspond to good coverage probabilities and  $Q^2$  coefficients. This can be seen in the case

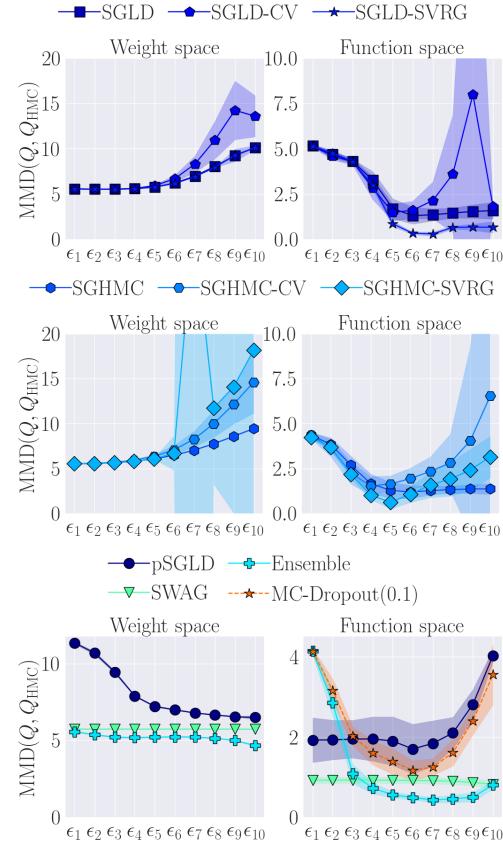


Figure 3. Problem AF#1. Maximum mean discrepancies  $MMD(\mathbb{Q}, \mathbb{Q}_{HMC})$  between the approximated posterior distributions  $\mathbb{Q}$  and the reference HMC sample, in weight and function spaces.

of **SWAG**, which achieves the lowest KSD values but poor coverage probabilities. Amongst the **SGMCMC** methods, it can also be seen that the lowest KSD values are often obtained for the lowest step sizes. As a result, the KSD seems often unable to identify approximation methods that generate valid confidence intervals. Several studies have shown that the KSD suffers from strong pathologies in sim-

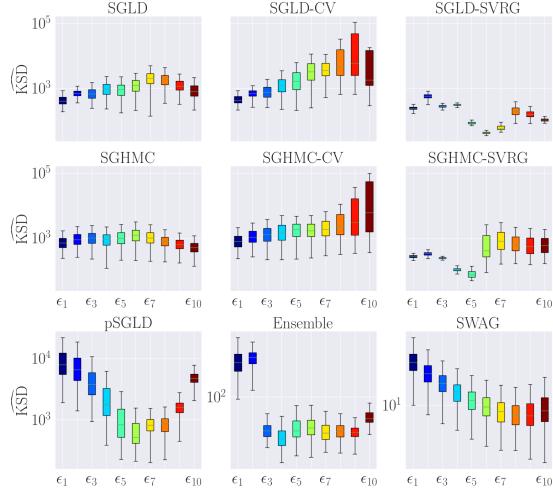


Figure 4. Problem AF#1. Kernelized Stein discrepancies  $KSD(\mathbb{P}, \mathbb{Q})$  between the target posterior measure  $\mathbb{P}$  and the approximated posteriors  $\mathbb{Q}$ . KSD values for the cyclical SGMCMC methods are reported in the Appendix.

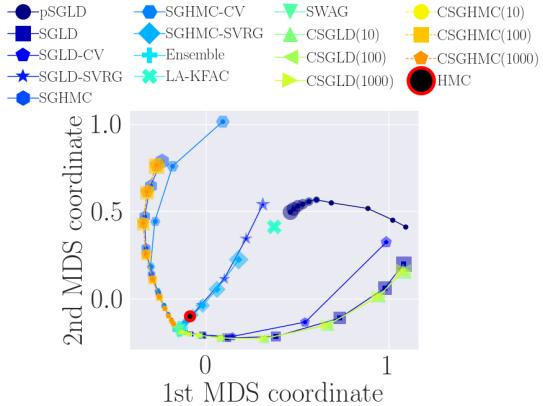


Figure 5. Dataset AF#1. Similarities between the algorithms as measured by the MMD in weight space and represented in a two-dimensional space build with multidimensional scaling. The markers sizes are proportional to the value of the step size  $\epsilon$ .

ple experiments (Wenliang & Kanagawa, 2020; Korba et al., 2021). In particular, the KSD is known to be insensitive to weight proportions in multimodal distributions, which may explain our observations.

A similar behavior is observed with the MMD distance to the HMC reference in weight space. For instance, this distance increases with the step size for **SGLD** but the best predictive performances are obtained for the highest step sizes. In contrast, the MMD distance to the HMC reference in function space is correlated to the predictive performance. A low MMD distance in function space typically corresponds to the best coverage probabilities and  $Q^2$  coefficients obtained by a given approximation method.

Finally, it can be observed that in this experiment, **deep ensembles**, **SGLD-SVRG**, and **SGHMC-SVRG** provide similar approximations as shown by the similarity measures in Figure 5, the MMD distances in function space, and the performances in terms of coverage probabilities and mean predictions. While the **LA-KFAC** method yields the highest KSD values and higher MMD distances than other methods, it offers a high precision on the test dataset, and a high coverage probability way above the target level.

## 7. Regression with OOD testing

In this section, we present the results obtained for the second regression problem (AF#2 in Table 1). Here, the distribution of test dataset differs from the distribution of the training dataset, which is more challenging in terms of coverage probabilities and predictive performances. We find that the behavior of the KSD with respect to the step size and approximation method is similar to the one observed in the previous experiment (see Figure 4). Furthermore, we also observe that the similarities between the approximation methods is close to the one observed in Figure 5. All the associated figures are reported in Appendix C for brevity.

**Coverage probabilities.** The coverage probabilities metrics and the regression coefficients  $Q^2$  are reported in Figure 6. Compared to the previous experiment of section 6, we see that achieving the target coverage level is more challenging for most of the approximation methods due to the test inputs that are out of the distribution of the training data. Several SGMCMC methods (such as **SGLD**, **SGHMC**, and their cyclical variants) easily yield high regression coefficients but struggle to reach the target confidence level. A similar behavior is observed for **deep ensembles**. This suggests that the width of the confidence intervals are not big enough to properly cover the test predictions. In this experiment, we see that **MC-dropout** has slightly better performances than in the previous section but remains less effective than the other approximation methods. Finally, we find that the performances of **pSGLD** and **SWAG** remain similar to AF#1.

**MMD distances to the HMC reference.** In terms of MMD distances, we find that **deep ensembles** and **SGLD-SVRG** yield the closest approximations to the HMC reference in function space. For brevity, the graphs of the MMD distances have been reported in Appendix C (see Figure 17). Here again, we find that the MMD distance in weight space increases with the step size while it mainly decreases in function space.

## 8. Discussion about coverage probabilities

Previous works that assess the quality of the confidence intervals in BNNs rely on the PICP (see, e.g., Yao et al.

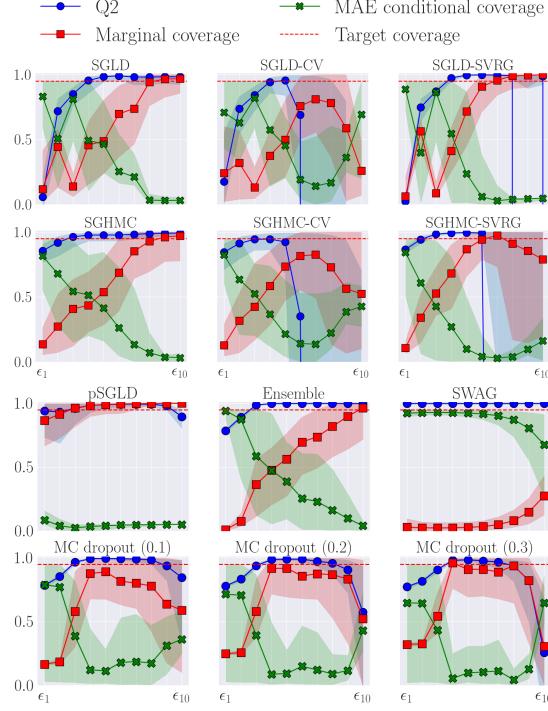


Figure 6. Problem AF#2. Coverage metrics and  $Q^2$  coefficient with respect to the step size  $\epsilon$ . The target coverage is set to 0.95.

(2019)). In contrast to MCP and CCP, PICP only integrates over the test set and does not take into account the variability of the confidence interval with respect to the training set. This is convenient in practice given that we generally do not have access to the distribution of the training data. However, PICP may lead to wrong conclusions about the validity of the confidence intervals. To support our claim, in Figure 7 we compare the three coverage metrics: the PICP histogram obtained by considering all the PICPs (mean over test set) for each training set, while the CPP histogram gathers all the CCPs (mean over training set) for each sample from the test set. We also represent the MCP (mean over both training and test set) and the target confidence level. In both cases, the MCP and CPP metrics clearly indicate that the target coverage is far from being attained. On the contrary, the PICP exhibits large variability over the training sets, with a significant probability of being close to the target coverage for **SGLD** on the left panel. This implies that, depending on the random sampling of the training set, the PICP may lead to consider **SGLD** as an efficient approximation of the posterior, while in reality the MCP and CPP are way more pessimistic and tend to conclude the opposite. Consequently, we advocate to use the PICP very cautiously.

Since in practice the MCP and CPP cannot be computed by sampling many independent training sets as we did in our experiments, we suggest instead to estimate them by resam-

pling methods such as the bootstrap, in order to compensate for the PICP limitations illustrated above.

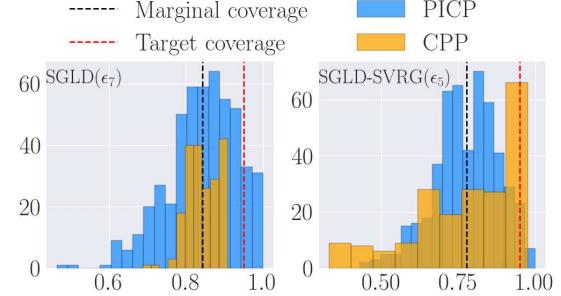


Figure 7. Problem AF#1. Comparison of the three coverage probabilities obtained with **SGLD** and **SGLD-SVRG**.

## 9. Summary

This work compares several approximation methods for Bayesian neural networks with 14 algorithms including 9 SGMC methods, **deep ensembles**, **LA-FKAC**, **MC-Dropout**, **SWAG**, and **HMC** which is taken as a reference. In contrast to previous works, we assess the performance of the selected methods by evaluating the validity of the generated confidence intervals, the quality of the approximation in weight and function spaces, and the precision of the mean prediction on test data. We find that **SGMCMC** methods and **deep ensembles** are able to achieve good coverage probabilities and predictive performances. Most notably, in our experiments, **deep ensembles** provide the best approximation to the HMC reference in both weight and function spaces. While **SWAG** has excellent performances in terms of regression, it tends to underestimate the uncertainties. We also propose a novel similarity analysis between the algorithms which shows that **LA-KFAC** and **pSGLD** with high step sizes provide close approximations in weight space. **SGMCMC-SVRG** and **deep ensembles** are also very similar in weight space. We finally consider the KSD but unfortunately we observe that it cannot exclusively be relied on to assess the quality of an approximation to the target posterior. In summary, our experiments show that several approximations methods can be efficient, but selecting appropriate hyperparameters yielding valid confidence intervals and regression accuracy is still challenging. A promising way would be to rely on resampling methods on the training set to estimate the MCP and the CCP or use the recently introduced conformal prediction (Barber et al., 2021). We also hope that the proposed new comparisons in weight and function spaces will give insights on how the different algorithms explore the posterior distribution.

## References

- Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., U.R., A., Makarenkov, V., and Nahavandi, S. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76:243–297, 2021.
- Ahn, S., Korattikara, A., and Welling, M. Bayesian posterior sampling via stochastic gradient fisher scoring. *arXiv preprint arXiv:1206.6380*, 2012.
- Baker, J., Fearnhead, P., Fox, E., and Nemeth, C. Control variates for stochastic gradient mcmc. *Statistics and Computing*, 29(3):599–615, 2019.
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. Predictive inference with the jackknife+. *The Annals of Statistics*, 49(1):486 – 507, 2021. doi: 10.1214/20-AOS1965. URL <https://doi.org/10.1214/20-AOS1965>.
- Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. Weight uncertainty in neural network. In *International conference on machine learning*, pp. 1613–1622. PMLR, 2015.
- Bradbury, J., Frostig, R., Hawkins, P., Johnson, M. J., Leary, C., Maclaurin, D., Necula, G., Paszke, A., VanderPlas, J., Wanderman-Milne, S., and Zhang, Q. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- Chen, T., Fox, E., and Guestrin, C. Stochastic gradient hamiltonian monte carlo. In *International conference on machine learning*, pp. 1683–1691, 2014.
- Cobb, A. and Jalaian, B. Scaling hamiltonian monte carlo inference for bayesian neural networks with symmetric splitting. In *Uncertainty in Artificial Intelligence*, pp. 675–685. PMLR, 2021.
- Daxberger, E., Kristiadi, A., Immer, A., Eschenhagen, R., Bauer, M., and Hennig, P. Laplace redux—effortless Bayesian deep learning. In *NeurIPS*, 2021.
- Ding, N., Fang, Y., Babbush, R., Chen, C., Skeel, R., and Neven, H. Bayesian sampling using stochastic gradient thermostats. *Advances in neural information processing systems*, 27, 2014.
- Dubey, K., J. Reddi, S., Williamson, S., Poczos, B., Smola, A., and Xing, E. Variance reduction in stochastic gradient langevin dynamics. *Advances in neural information processing systems*, 29:1154–1162, 2016.
- Fort, S., Hu, H., and Lakshminarayanan, B. Deep ensembles: A loss landscape perspective. *arXiv preprint arXiv:1912.02757*, 2019.
- Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059. PMLR, 2016.
- Gawlikowski, J., Tassi, C., Ali, M., Lee, J., Humt, M., Feng, J., Kruspe, A., Triebel, R., Jung, P., Roscher, R., et al. A survey of uncertainty in deep neural networks. *arXiv preprint arXiv:2107.03342*, 2021.
- Girolami, M. and Calderhead, B. Riemann manifold langevin and hamiltonian monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(2):123–214, 2011.
- Gorham, J. and Mackey, L. Measuring sample quality with kernels. In *International Conference on Machine Learning*, pp. 1292–1301. PMLR, 2017.
- Graves, A. Practical variational inference for neural networks. *Advances in neural information processing systems*, 24, 2011.
- Gretton, A., Borgwardt, K., Rasch, M., Schölkopf, B., and Smola, A. A kernel method for the two-sample-problem. *Advances in neural information processing systems*, 19, 2006.
- Hoffman, M., Blei, D., Wang, C., and Paisley, J. Stochastic variational inference. *Journal of Machine Learning Research*, 2013.
- Izmailov, P., Vikram, S., Hoffman, M., and Wilson, A. What are bayesian neural network posteriors really like? In *International Conference on Machine Learning*, pp. 4629–4640. PMLR, 2021.
- Korba, A., Aubin-Frankowski, P.-C., Majewski, S., and Ablin, P. Kernel stein discrepancy descent. In *International Conference on Machine Learning*, pp. 5719–5730. PMLR, 2021.
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- Lao, J. and Louf, R. Blackjax: A sampling library for JAX, 2020. URL <http://github.com/blackjax-devs/blackjax>.
- Li, C., Chen, C., Carlson, D., and Carin, L. Preconditioned stochastic gradient langevin dynamics for deep neural networks. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.

- Lin, Z., Trivedi, S., and Sun, J. Locally valid and discriminative prediction intervals for deep learning models. *Advances in Neural Information Processing Systems*, 34: 8378–8391, 2021.
- Liu, Q., Lee, J., and Jordan, M. A kernelized stein discrepancy for goodness-of-fit tests. In *International conference on machine learning*, pp. 276–284. PMLR, 2016.
- Louizos, C. and Welling, M. Multiplicative normalizing flows for variational bayesian neural networks. In *International Conference on Machine Learning*, pp. 2218–2227. PMLR, 2017.
- Ma, Y., Chen, T., and Fox, E. A complete recipe for stochastic gradient mcmc. *Advances in neural information processing systems*, 28, 2015.
- Maddox, W., Izmailov, P., Garipov, T., Vetrov, D., and Wilson, A. A simple baseline for bayesian uncertainty in deep learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Neal, R. e. a. Mcmc using hamiltonian dynamics. *Handbook of Markov chain Monte Carlo*, 2(11):2, 2011.
- Nemeth, C. and Fearnhead, P. Stochastic gradient markov chain monte carlo. *Journal of the American Statistical Association*, 116(533):433–450, 2021.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.
- Ritter, H., Botev, A., and Barber, D. A scalable laplace approximation for neural networks. In *6th International Conference on Learning Representations, ICLR 2018-Conference Track Proceedings*, volume 6. International Conference on Representation Learning, 2018.
- Sejdinovic, D., Sriperumbudur, B., Gretton, A., and Fukumizu, K. Equivalence of distance-based and rkhs-based statistics in hypothesis testing. *The Annals of Statistics*, pp. 2263–2291, 2013.
- Springenberg, J., Klein, A., Falkner, S., and Hutter, F. Bayesian optimization with robust bayesian neural networks. *Advances in neural information processing systems*, 29:4134–4142, 2016.
- Teymur, O., Gorham, J., Riabiz, M., and Oates, C. Optimal quantisation of probability measures using maximum mean discrepancy. In *International Conference on Artificial Intelligence and Statistics*, pp. 1027–1035. PMLR, 2021.
- Torgerson, W. S. Multidimensional scaling: I. theory and method. *Psychometrika*, 17(4):401–419, 1952.
- Welling, M. and Teh, Y. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688. Citeseer, 2011.
- Wenliang, L. and Kanagawa, H. Blindness of score-based methods to isolated components and mixing proportions. *arXiv preprint arXiv:2008.10087*, 2020.
- Wenzel, F., Roth, K., Veeling, B., Swiatkowski, J., Tran, L., Mandt, S., Snoek, J., Salimans, T., Jenatton, R., and Nowozin, S. How good is the Bayes posterior in deep neural networks really? In III, H. D. and Singh, A. (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 10248–10259. PMLR, 13–18 Jul 2020.
- Yao, J., Pan, W., Ghosh, S., and Doshi-Velez, F. Quality of uncertainty quantification for bayesian neural network inference. *arXiv preprint arXiv:1906.09686*, 2019.
- Zhang, R., Li, C., Zhang, J., Chen, C., and Wilson, A. Cyclical stochastic gradient mcmc for bayesian deep learning. *arXiv preprint arXiv:1902.03932*, 2019.

# Supplementary Material

## A. Considered algorithms

We briefly synthesize the considered **SGMCMC** methods that we consider in our study, with details on software implementation and practical considerations in our numerical experiments.

### A.1. Stochastic gradient Markov chain Monte Carlo algorithms

Stochastic gradient Markov chain Monte Carlo (SGMCMC) methods are a family of gradient-based algorithms that rely on the resolution of an Itô stochastic differential equation (Ma et al., 2015) (see also (Nemeth & Fearnhead, 2021) for a review). Amongst this family of methods, we consider the stochastic gradient Langevin dynamics (SGLD) (Welling & Teh, 2011) and the stochastic gradient Hamiltonian Monte Carlo (SGHMC) (Chen et al., 2014) methods, together with other variants that include variance reduction (Dubey et al., 2016; Baker et al., 2019) or preconditioning techniques (Li et al., 2016; Springenberg et al., 2016).

**Stochastic gradient Langevin dynamics.** The **SGLD** algorithm proposed by Welling & Teh (2011) takes the form

$$k \geq 0, \quad \mathbf{w}^{k+1} = \mathbf{w}^k - \epsilon_k \hat{\nabla} U(\mathbf{w}^k) + \sqrt{2\epsilon_k} \Delta \mathbf{W}^{k+1},$$

where  $\Delta \mathbf{W}^{k+1}$  is centered and normalized Gaussian random variable. Herein, **SGLD** is implemented with a constant step size  $\epsilon_k = \epsilon$  for all  $k \geq 0$ .

**Stochastic gradient Hamiltonian Monte Carlo.** The **SGHMC** algorithm has been proposed by Chen et al. (2014) where one leapfrog step reads as:

$$k \geq 0 : \begin{cases} \mathbf{w}^{k+1} = \mathbf{w}^k + \eta_k \mathbf{M}^{-1} \mathbf{r}^k, \\ \mathbf{r}^{k+1} = \mathbf{r}^k - \eta_k \hat{\nabla} U(\mathbf{w}^k) - \epsilon_k \mathbf{C} \mathbf{M}^{-1} \mathbf{r}^k + \sqrt{2\eta_k} \Delta \mathbf{W}^{k+1}, \end{cases}$$

where  $\mathbf{r}$  denotes the momentum and  $\Delta \mathbf{W}^{k+1} \sim \mathcal{N}(0, \mathbf{I})$ . By introducing the change of variable  $\mathbf{v}^k = \epsilon_k \mathbf{M}^{-1} \mathbf{r}^k$ , the above equations can be rewritten as

$$k \geq 0 : \begin{cases} \mathbf{w}^{k+1} = \mathbf{w}^k + \mathbf{v}^k, \\ \mathbf{v}^{k+1} = (1 - \alpha) \mathbf{v}^k - \epsilon_k \hat{\nabla} U(\mathbf{w}^k) + \sqrt{2\alpha\epsilon_k} \Delta \mathbf{W}^{k+1}, \end{cases}$$

where  $\epsilon_k = \eta_k^2 \mathbf{M}^{-1}$  and  $\alpha = \eta \mathbf{M}^{-1} \mathbf{C}$ . The step size  $\epsilon_k$  is constant in our experiments with **SGHMC**, and the momentum  $\mathbf{v}$  is resampled every 10 leapfrog steps following a centered and normalized Gaussian distribution.

**Preconditioned SGLD.** (Ma et al., 2015) proposed the Riemannian SGLD by building upon the Riemannian manifold Langevin and Hamiltonian Monte Carlo methods (Girolami & Calderhead, 2011). The Riemannian SGLD takes the form:

$$k \geq 0, \quad \mathbf{w}^{k+1} = \mathbf{w}^k - \epsilon_k \left( D(\mathbf{w}^k) \hat{\nabla} U(\mathbf{w}^k) + \boldsymbol{\Gamma}(\mathbf{w}^k) \right) + \sqrt{2\epsilon_k D(\mathbf{w}^k)} \Delta \mathbf{W}^{k+1}, \quad (1)$$

where  $\mathbf{D} \in \mathbb{R}^{d \times d}$  can be seen as a preconditioning matrix and the vector  $\boldsymbol{\Gamma} \in \mathbb{R}^d$  is given by

$$\boldsymbol{\Gamma}_i(\mathbf{w}) = \sum_j \frac{\partial}{\partial \zeta_j} \mathbf{D}_{ij}(\mathbf{w}). \quad (2)$$

Several forms of the preconditioner  $\mathbf{D}$  can be used, the most common choice being a constant diagonal matrix  $\mathbf{D}_{\text{diag}}$  of the form  $\mathbf{D}_{\text{diag}} = N^{-1} \mathbf{I}_d$ , where  $N$  denotes the size of the training dataset. The expected Fisher information matrix is also a natural choice but it is usually intractable, and a few alternatives have been proposed in the literature. Herein we consider

the **pSGLD** method of Li et al. (2016) that neglects the vector  $\Gamma$  and uses a preconditioner inspired by the RMSprop optimization algorithm. The gradient is scaled using a moving average of its norm at each iteration:

$$\mathbf{D}(\mathbf{w}^k) = \text{diag} \left( \lambda \mathbf{I} + \sqrt{\mathbf{L}(\mathbf{w}^k)} \right)^{-1},$$

$$\mathbf{L}(\mathbf{w}^k) = \alpha \mathbf{L}(\mathbf{w}^{k-1}) + (1 - \alpha) \widehat{\nabla} U(\mathbf{w}^{k-1}) \circ \widehat{\nabla} U(\mathbf{w}^{k-1}),$$

where  $\alpha$  is a parameter with values in  $[0, 1]$ ,  $\lambda$  is a regularization constant, and  $\circ$  denotes the Hadamard (element-wise) product.

**SGMCMC-CV and SGMCMC-SVRG.** As described in section 3, we consider two variants of the **SGLD** and **SGHMC** algorithms that rely on control variates to reduce the variance of stochastic approximation  $\widehat{\nabla} U$  of  $\nabla U$ . Following Baker et al. (2019), **SGLD-CV** and **SGHMC-CV** use an estimation of the maximum a posteriori (MAP) for the control variable  $\eta$ . Hence, the neural network is first trained in a classical fashion using an optimization algorithm, and the resulting values of the network parameters are stored. The remaining variants **SGLD-SVRG** and **SGHMC-SVRG** use the strategy proposed by Dubey et al. (2016) which consists in setting  $\eta$  to a MAP estimate as well, but then updating  $\eta$  every  $m$  iterations as follows:  $\eta = \mathbf{w}^\ell$  if  $\text{mod}(\ell, m) = 0$ . The update frequency  $m$  has been fixed to 100 in every experiment.

**Cyclical SGMCMC.** The cyclical variants of **SGLD** and **SGHMC** use a step size of the form (Zhang et al., 2019)

$$\epsilon_k = \frac{\epsilon_0}{2} \left( \cos \left( \frac{\pi \text{mod}(k-1, \lceil K/M \rceil)}{\lceil K/M \rceil} \right) + 1 \right),$$

where  $\epsilon_0$  is the initial step size,  $K$  denotes total number of iterations, and  $M$  is the cycle length. The performance of the cyclical variants is studied with respect to the initial step size and the cycle length, as described in section 4.

## A.2. Implementation

Our experiments are implemented with JAX (Bradbury et al., 2018) and PyTorch (Paszke et al., 2019). For the **LA-FKAC** approximation, we use the Laplace package that implements various types of Laplace approximations (Daxberger et al., 2021) and relies on PyTorch. The remaining methods are implemented with Google’s JAX framework. Regarding the **SGMCMC** methods, we rely on BlackJAX (Lao & Louf, 2020), where **SGLD** and **SGHMC** are already implemented. We use our own implementations of variants with variance reduction, integrated within BlackJAX. The remaining approximation methods are implemented from scratch with JAX.

## A.3. Hyperparameters for numerical experiments

The selected step sizes for each **SGMCMC** algorithm are given in Table 3. Any other hyperparameter, such as the number of iterations or batch size, is held fixed. For all **SGMCMC** methods, the number of iterations is set to  $10^5$  and the first  $50^4$  iterations are discarded as burn-in. The remaining iterations are automatically thinned by selecting 2000 particles which are obtained by minimizing the maximum mean discrepancy (Gretton et al., 2006). This procedure, referred to as MMD thinning, is applied to every **SGMCMC** output and works as follows. Given  $T$  samples  $\mathbf{w}_1, \dots, \mathbf{w}_T$  of a network parameters  $\mathbf{w} \in \mathbb{R}^d$ , we find a sequence of indices  $\pi \in \{1, \dots, T\}^m$  such that  $\mathbf{w}_{\pi(1)}, \dots, \mathbf{w}_{\pi(m)}$  represent the selected samples. The sequence of indices  $\pi$  is obtained thanks to a greedy quantization algorithm (Teymur et al., 2021) that minimizes the maximum mean discrepancy. Let then  $\mathbb{P}_T$  be the empirical distribution of the  $T$  samples  $\mathbf{w}_1, \dots, \mathbf{w}_T$ ,  $\mathbb{P}_T = \sum_{i=1}^T \delta(\mathbf{w}_i)$ . At the  $(i+1)$ -th iteration of the quantization algorithm, the index  $\pi(i+1)$  is obtained by minimizing the MMD as follows:

$$\pi(i+1) = \arg \min_{j \in \{1, \dots, T\}} \text{MMD}^2(\mathbb{P}_T, \mathbb{Q}_{i+1}(j)),$$

where  $\mathbb{Q}_{i+1}(j)$  denotes the empirical measure of the already  $i$  selected samples  $\mathbf{w}_1, \dots, \mathbf{w}_i$ , and an additional sample  $\mathbf{w}_j$  to be determined:

$$\mathbb{Q}_{i+1}(j) = \frac{1}{i+1} \sum_{\ell=1}^i \delta(\mathbf{w}_{\pi(\ell)}) + \frac{1}{i+1} \delta(\mathbf{w}_j).$$

The underlying kernel function  $k$  is chosen as the following characteristic distance-based kernel (Sejdinovic et al., 2013)  $k(\mathbf{w}, \mathbf{w}') = \|\mathbf{w}\|_2 + \|\mathbf{w}'\|_2 - \|\mathbf{w} - \mathbf{w}'\|_2$ . In our all experiments, we select  $m = 2000$  samples in order to represent each

**SGMCMC** output. In this case of the variance reduction technique **SVRG**, the control variates are updated every 100 iterations. We use 10 leapfrog steps in **SGHMC** and each of its variants. For **deep ensembles**, 200 neural networks are independently trained and their weights are initialized from a centered normalized Gaussian distribution. For each of the remaining approximation methods (**LA-KFAC**, **MC-Dropout**, and **SWAG**), samples of sizes 2000 are saved.

Table 3. Selected step sizes for each algorithm, equally log-spaced.

Step size	SGLD variants	SGHMC variants	pSGLD, Ensemble, MC Dropout	LA-KFAC	SWAG
$\epsilon_1$	1e-08	1e-08	0.0001	0.000001	1.000000e-07
$\epsilon_2$	2.154435e-08	1.668101e-08	0.000215	0.000004	1.668101e-07
$\epsilon_3$	4.641589e-08	2.782559e-08	0.000464	0.000013	2.782559e-07
$\epsilon_4$	1.000000e-07	4.641589e-08	0.001000	0.000046	4.641589e-07
$\epsilon_5$	2.154435e-07	7.742637e-08	0.002154	0.000167	7.742637e-07
$\epsilon_6$	4.641589e-07	1.291550e-07	0.004642	0.000599	1.291550e-06
$\epsilon_7$	1.000000e-06	2.154435e-07	0.010000	0.002154	2.154435e-06
$\epsilon_8$	2.154435e-06	3.593814e-07	0.021544	0.007743	3.593814e-06
$\epsilon_9$	4.641589e-06	5.994843e-07	0.046416	0.027826	5.994843e-06
$\epsilon_{10}$	1e-05	1e-06	0.1	0.1	1e-05

## B. Description of the training and test datasets

The synthetic regression problems are described in more details below.

**Regression problem AF#1.** This first regression test case is a homoscedastic regression of a one-dimensional function. The dataset  $\mathcal{D} = \{(X_i, Y_i)\}_{i=1}^N$  is such that  $X_i \sim \mathcal{U}([-3, 3])$ , and  $Y_i = \cos(2X_i) + \sin(X_i) + \epsilon_i$ , where  $\epsilon_i \sim \mathcal{N}(0, \sigma)$ , with  $\sigma = 0.2$ . The test dataset  $\mathcal{D}^*$  has the same distribution.

**Regression problem AF#2.** This test case is taken from (Yao et al., 2019) which consists in a homoscedastic regression of a one-dimensional function as well. The inputs  $X_1, \dots, X_N$  of a training dataset  $\mathcal{D}$  are uniformly sampled from  $[-4, -1] \cup [1, 4]$ , while test inputs are uniformly distributed between  $[-4, 4]$ . The observations are defined as  $Y_i = 0.1X_i^3 + \epsilon_i$  where  $\epsilon_i \sim \mathcal{N}(0, 0.25)$ .

**Regression problem AF#3.** This test case is taken from (Yao et al., 2019) which consists in a homoscedastic regression of a one-dimensional function. The first 80 inputs  $X_1, \dots, X_{80}$  of a training dataset  $\mathcal{D}$  are uniformly sampled in  $[-6, -2] \cup [2, 6]$ , and the last two samples  $X_{81}, X_{82}$  are uniformly sampled in  $[-2, 2]$ . In contrast, the inputs of a test dataset  $\mathcal{D}^*$  are all uniformly distributed in  $[-6, 6]$ . The output observations are given by  $Y_i = -(1+X_i) \sin(1.2X_i) + \epsilon_i$  where  $\epsilon \sim \mathcal{N}(0, 0.25)$ .

**Regression problem AF#4.** This last regression problem is taken from the work of Izmailov et al. (2021). The inputs of a training dataset  $\mathcal{D} = \{(X_i, Y_i)\}_{i=1}^N$  are uniformly distributed in  $[-10, -6] \cup [6, 10] \cup [14, 18]$ . The observations  $Y_i$  are defined as  $Y_i = f(X_i; \mathbf{w}_0) + \epsilon_i$ , with  $\epsilon_i \sim \mathcal{N}(0, 0.02)$ , and where  $f(\cdot; \mathbf{w}_0)$  denotes a feed-forward neural network with three hidden layers of sizes 100 (and thus 20,501 parameters in total). The weight parameters  $\mathbf{w}_0$  are sampled once from a centered normalized Gaussian distribution  $\mathcal{N}(0, \mathbf{I}_d)$ . The inputs of the test dataset are uniformly distributed in  $[-12, 22]$ .

## C. Additional results

In this section, we gather additional results about the empirical experiments, and a few additional details about the networks architecture. For the first synthetic problem (AF #1), the architecture is made of two hidden layers with 100 hidden features, leading to a network with  $d = 10,401$  parameters. Following (Yao et al., 2019), we consider two hidden layers with 50 features for the second and third regression problems (AF #2 and AF #3), such that the network is made of  $d = 2,651$  parameters. Finally, following Izmailov et al. (2021), a network made of three hidden layers with 100 features each is used for the last synthetic problem (AF #4), leading to a network made of  $d = 20,501$  parameters.

### C.1. Best marginal coverage probabilities and regression coefficients

We first gather below in Tables 4-6 the best marginal coverages and best  $Q^2$  coefficients obtained for the experiments AF#2, AF#3, and AF#4.

Table 4. AF#2: Best marginal coverage probabilities (MCP) and best  $Q^2$  coefficients obtained with each algorithm.

Method	Best MCP	$Q^2$	Best $Q^2$ MCP	Method	Best MCP	$Q^2$	Best $Q^2$ MCP
SGLD	<b>0.94</b>	0.98	0.989 0.696	CSGHMC(10)	<b>0.947</b>	0.988	0.988 0.947
SGLD-CV	0.809	-4e8	0.954 0.497	CSGHMC(100)	<b>0.946</b>	0.990	<b>0.99</b> 0.946
SGLD-SVRG	<b>0.954</b>	0.995	<b>0.997</b> 0.908	CSGHMC(1000)	0.928	0.99	<b>0.99</b> 0.928
SGHMC	<b>0.96</b>	0.99	<b>0.99</b> 0.960	pSGLD	0.962	0.961	<b>0.993</b> 0.996
SGHMC-CV	0.825	-13.05	0.944 0.583	Deep ensemble	0.964	0.997	<b>0.997</b> 0.736
SGHMC-SVRG	<b>0.942</b>	0.992	<b>0.995</b> 0.867	SWAG	0.273	0.997	<b>0.997</b> 0.273
CSGLD(10)	<b>0.949</b>	0.978	0.989 0.655	MC Drop.(0.1)	0.890	0.994	<b>0.994</b> 0.89
CSGLD(100)	0.92	0.976	0.99 0.652	MC Drop.(0.2)	0.917	0.989	0.989 0.917
CSGLD(1000)	0.91	0.980	0.989 0.648	MC Drop.(0.3)	<b>0.959</b>	0.979	0.984 0.909

Table 5. AF#3: Best marginal coverage probabilities (MCP) and best  $Q^2$  coefficients obtained with each algorithm.

Method	Best MCP	$Q^2$	Best $Q^2$ MCP	Method	Best MCP	$Q^2$	Best $Q^2$ MCP
SGLD	0.921	0.803	0.823 0.885	CSGHMC(10)	0.834	0.825	0.867 0.781
SGLD-CV	0.551	-29.357	-0.091 0.034	CSGHMC(100)	0.803	0.794	0.850 0.753
SGLD-SVRG	0.919	0.859	0.879 0.877	CSGHMC(1000)	0.805	0.817	0.851 0.740
SGHMC	0.837	0.788	0.847 0.801	pSGLD	0.919	0.851	0.879 0.907
SGHMC-CV	0.387	-28.251	-0.165 0.131	Deep ensemble	<b>0.94</b>	0.933	0.944 0.821
SGHMC-SVRG	0.736	-1e9	0.655 0.614	SWAG	0.751	0.538	nan 0.654
CSGLD(10)	0.902	0.829	0.833 0.856	MC Drop.(0.1)	0.624	0.792	0.799 0.588
CSGLD(100)	0.893	0.739	0.803 0.840	MC Drop.(0.2)	0.554	0.662	0.662 0.554
CSGLD(1000)	0.900	0.818	0.859 0.851	MC Drop.(0.3)	0.415	0.436	0.437 0.397

Table 6. AF#4: Best marginal coverage probabilities (MCP) and best  $Q^2$  coefficients obtained with each algorithm.

Method	Best MCP	$Q^2$	Best $Q^2$ MCP	Method	Best MCP	$Q^2$	Best $Q^2$ MCP
SGLD	0.646	-2.136	0.785 0.188	CSGHMC(10)	0.589	-2.316	0.814 0.133
SGLD-CV	0.696	-57.602	0.781 0.130	CSGHMC(100)	0.836	0.446	0.780 0.123
SGLD-SVRG	0.529	0.488	0.869 0.390	CSGHMC(1000)	0.475	0.189	0.830 0.266
SGHMC	0.661	-1.575	0.792 0.098	pSGLD	0.957	-20858.904	-324.041 1.0
SGHMC-CV	0.616	-1284.229	0.786 0.148	Deep ensemble	0.413	0.420	0.925 0.203
SGHMC-SVRG	0.607	-0.367	0.829 0.176	swag	0.537	0.862	0.862 0.537
CSGLD(10)	0.639	-40.265	0.769 0.087	MC Drop.(0.1)	0.232	0.802	0.851 0.227
CSGLD(100)	0.654	-5.929760	0.785 0.139	MC Drop.(0.2)	0.400	0.796	0.860 0.270
CSGLD(1000)	0.563	-0.639	0.769 0.126	MC Drop.(0.3)	0.453	0.549	0.821 0.314

### C.2. Similarities between the algorithms

The similarities between the algorithms are gathered for all the experiments in Figure 8. The similarities in weight space seem to exhibit the same structure in our four experiments. However, the similarities in function space do not exhibit any particular structure.

### C.3. Additional results for regression problem AF#1

Results obtained for the cyclical variants of **SGMCMC** have been reported herein. Figure 9 shows the  $Q^2$  regression coefficient and coverage probabilities obtained with **CSGMCMC**.

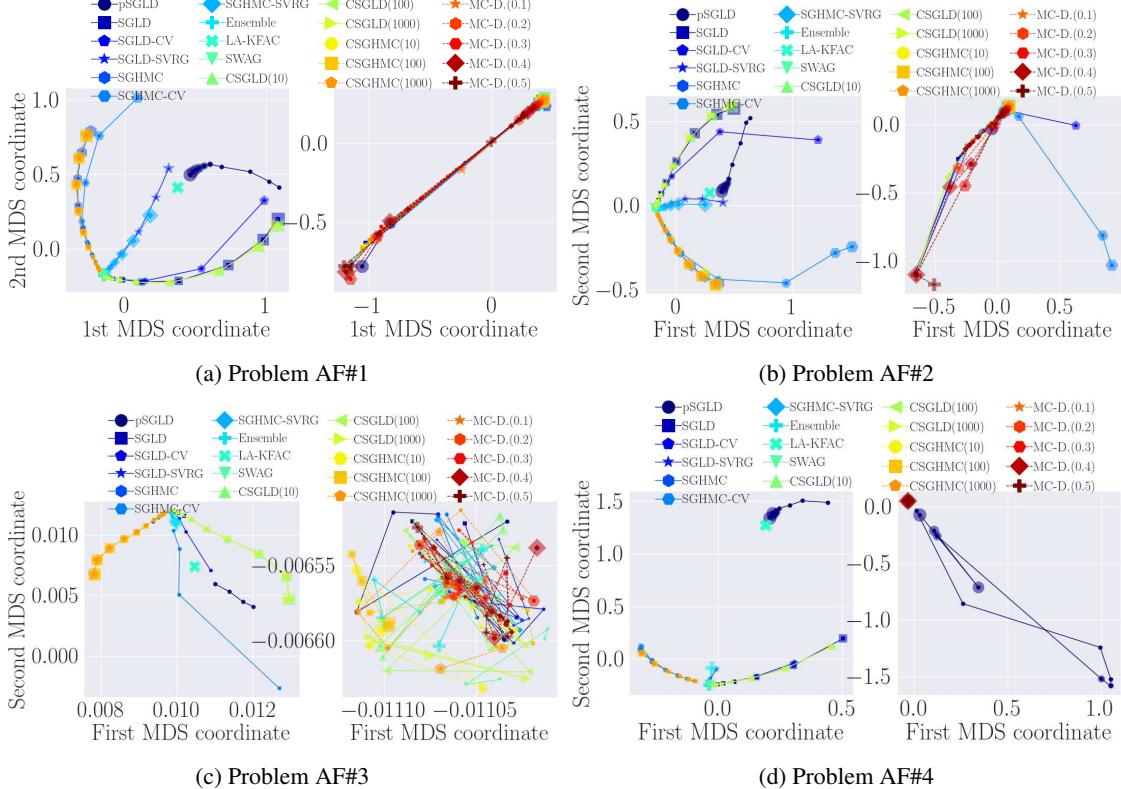


Figure 8. Similarities between the algorithms as measured by the MMD and represented in a two-dimensional space build with multi-dimensional scaling. Left panels: weight space, right panels: function space. The markers sizes are proportional to the value of the underlying step size  $\epsilon$ .

We also report here graphs of the marginal coverage probability with respect to the target confidence level  $(1 - \alpha)$ , for some  $\alpha \in [0.05, 0.95]$ , and all the considered values for the underlying hyperparameters, which are gathered in Figures 10-13.

We find that most **SGMCMC** methods (without variance reduction) give rise to marginal coverage probabilities that do not increase monotonically with the target level, and that are easily much higher than the target level. The cyclical step size reduces the marginal coverage probabilities. The **pSGLD** methods yield very different results as the target level is easily achieved, even for the lowest step sizes. **LA-KFAC** shows a similar behavior, where the marginal coverage probability is always above the target level.

The **MC-Dropout** method is able to reach the target level whenever  $\alpha$  is small enough, but struggles to reach an appropriate marginal coverage as  $\alpha$  increases. Finally, we find that the marginal coverage obtained with **deep ensembles** varies monotonically with the target level for sufficiently high step sizes.

#### C.4. Additional results for regression problem AF#2

The coverage metrics and regression coefficient obtained with **CSGLD** and **CSGHMC** are shown in Figure 13. The performances of the cyclical variants are similar to those of **SGLD** and **SGHMC** (see Figure 6). The graphs of the kernelized Stein discrepancies and maximum mean discrepancies are shown in Figures 14 and 17-18.

#### C.5. Additional results for regression problems AF#3 and AF#4

All the results obtained for the third and fourth regression problems are shown in Figures 19-24, and Figures 25-30. The distributions of the test data differ significantly from the training data in problem AF#3, and even more in problem AF#4. It can be seen in the associated results that these problems are challenging for all the methods, which are not able to reach the target coverage probability, nor obtain good regression coefficients.

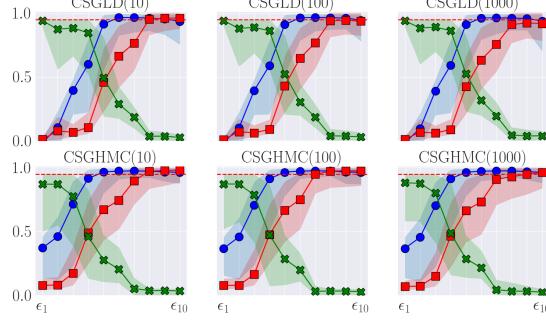


Figure 9. Problem AF#1. Coverage metrics and  $Q^2$  coefficient with respect to the step size  $\epsilon$ . The target coverage is set to 0.95.

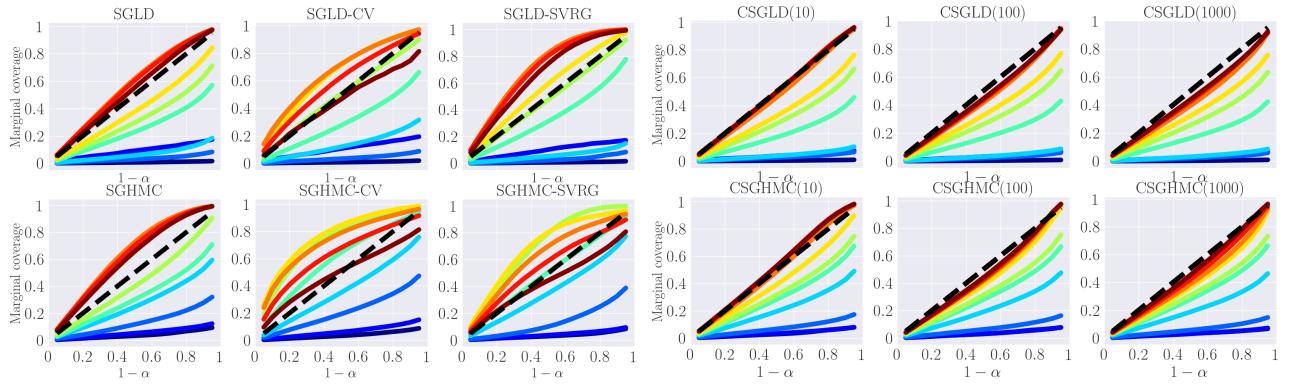


Figure 10. Problem AF#1. Graphs of the marginal coverage probability with respect to the target level for SGMCMC methods. The curves are colored by the value of the underlying step size: dark blue corresponds to the lowest, while dark red corresponds to the highest step size.

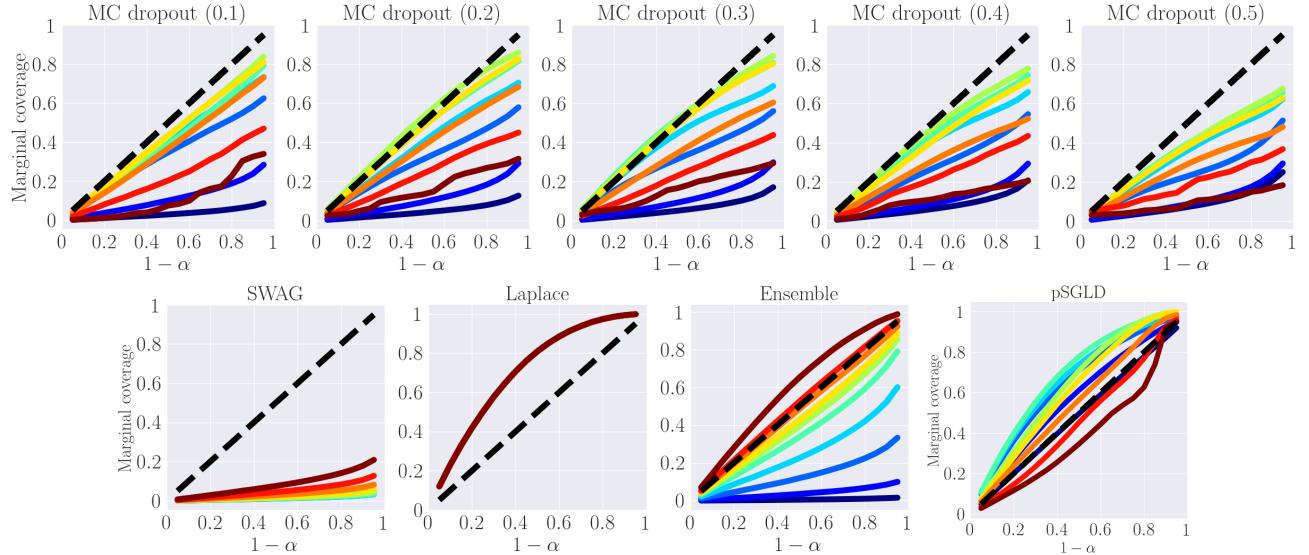


Figure 11. Problem AF#1. Graphs of the marginal coverage probability with respect to the target level for MC-Dropout, SWAG, LA-KFAC, deep ensembles, and pSGLD.

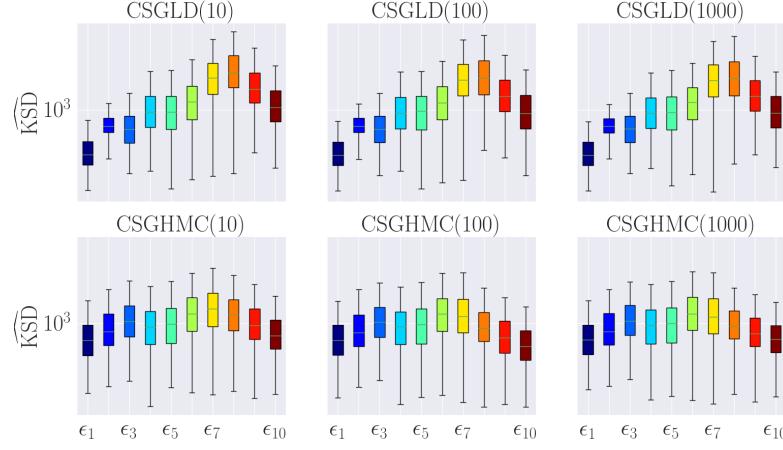


Figure 12. Problem AF#1. Kernelized Stein discrepancies  $\text{KSD}(\mathbb{P}, \mathbb{Q})$  between the target posterior measure  $\mathbb{P}$  and the approximated posteriors  $\mathbb{Q}$ .

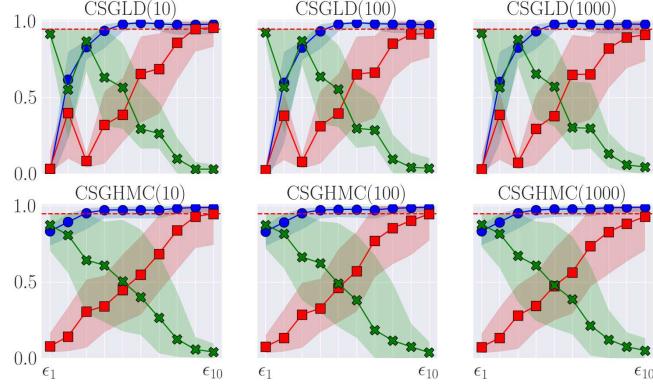


Figure 13. Problem AF#2. Coverage metrics and  $Q^2$  coefficient with respect to the step size  $\epsilon$ . The target coverage is set to 0.95.

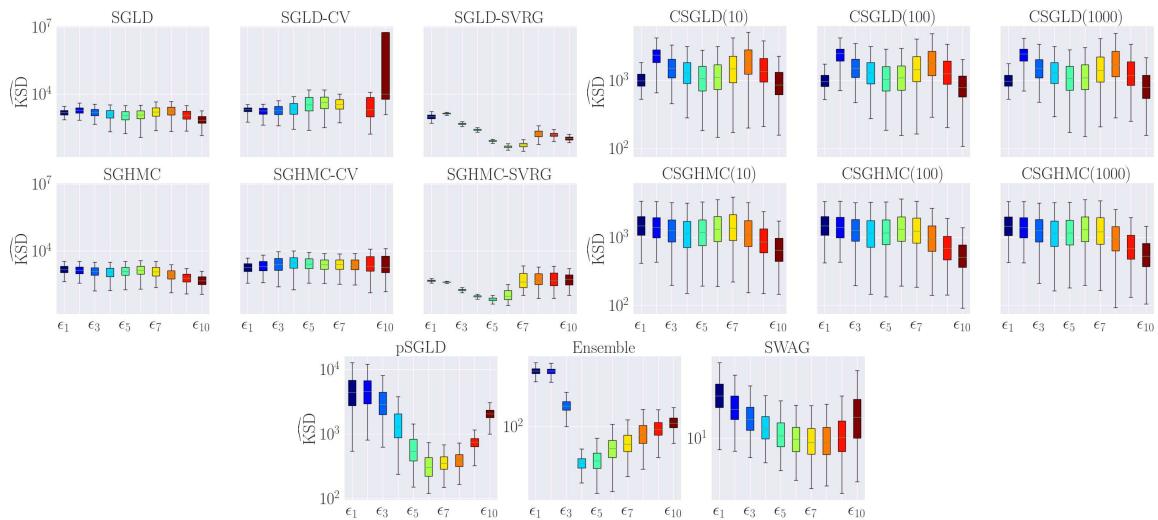


Figure 14. Problem AF#2. Kernelized Stein discrepancies  $\text{KSD}(\mathbb{P}, \mathbb{Q})$  between the target posterior measure  $\mathbb{P}$  and the approximated posteriors  $\mathbb{Q}$ .

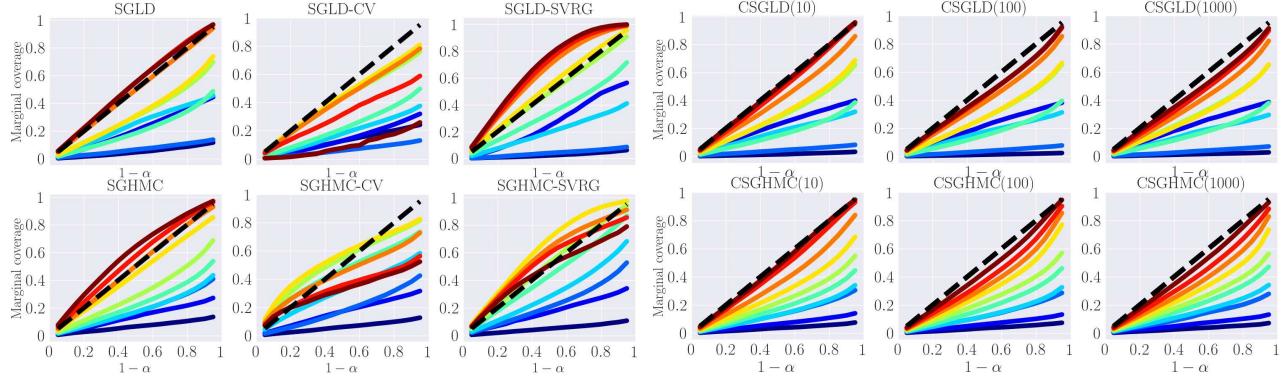


Figure 15. Problem AF#2. Graphs of the marginal coverage probability with respect to the target level for **SGMCMC** methods. The curves are colored by the value of the underlying step size: dark blue corresponds to the lowest, while dark red corresponds to the highest step size.

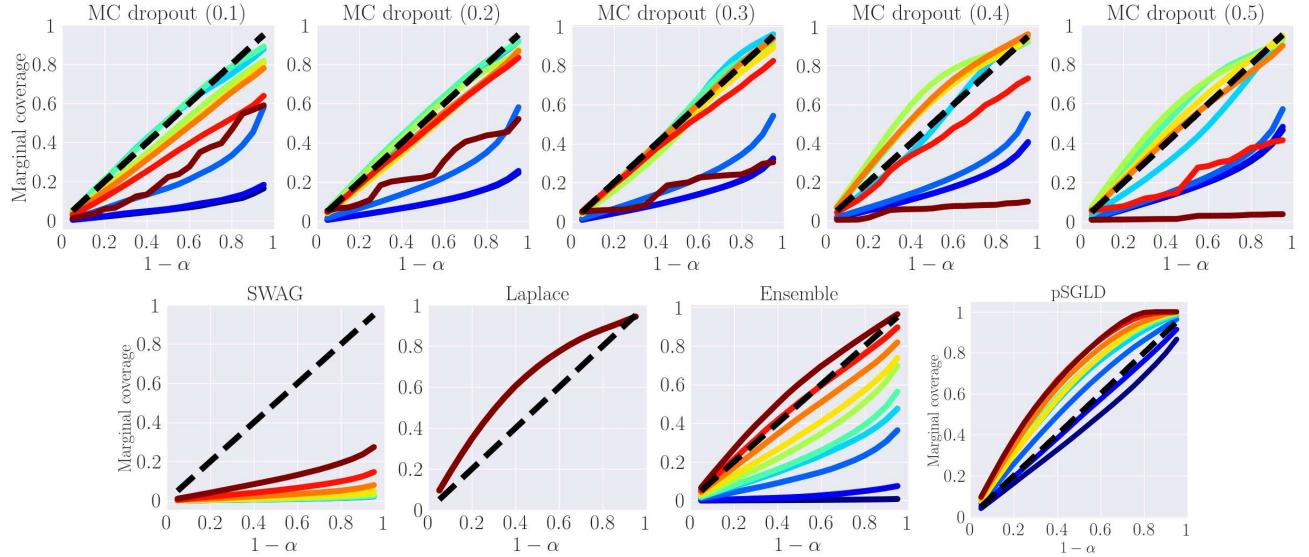


Figure 16. Problem AF#2. Graphs of the marginal coverage probability with respect to the target level for **MC-Dropout**, **SWAG**, **LA-KFAC**, **deep ensembles**, and **pSGLD**.

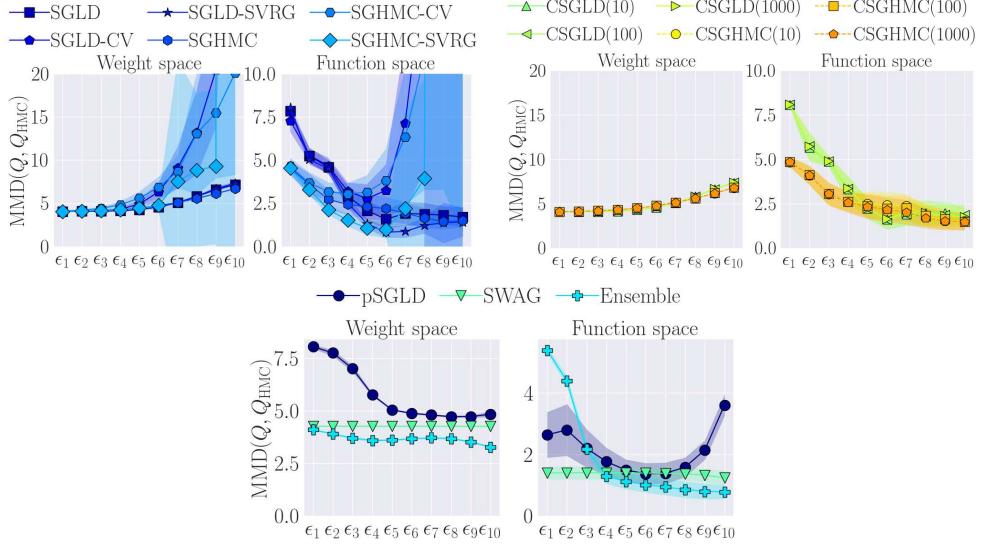


Figure 17. Problem AF#2. Maximum mean discrepancies  $MMD(\mathbb{Q}, \mathbb{Q}_{HMC})$  between the approximated posterior distributions  $\mathbb{Q}$  and the reference HMC sample, in weight and function spaces.

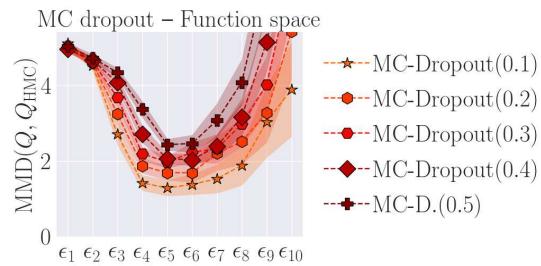


Figure 18. Problem AF#2. Maximum mean discrepancies between the approximated posterior distributions and the reference HMC sample, in weight and function spaces.

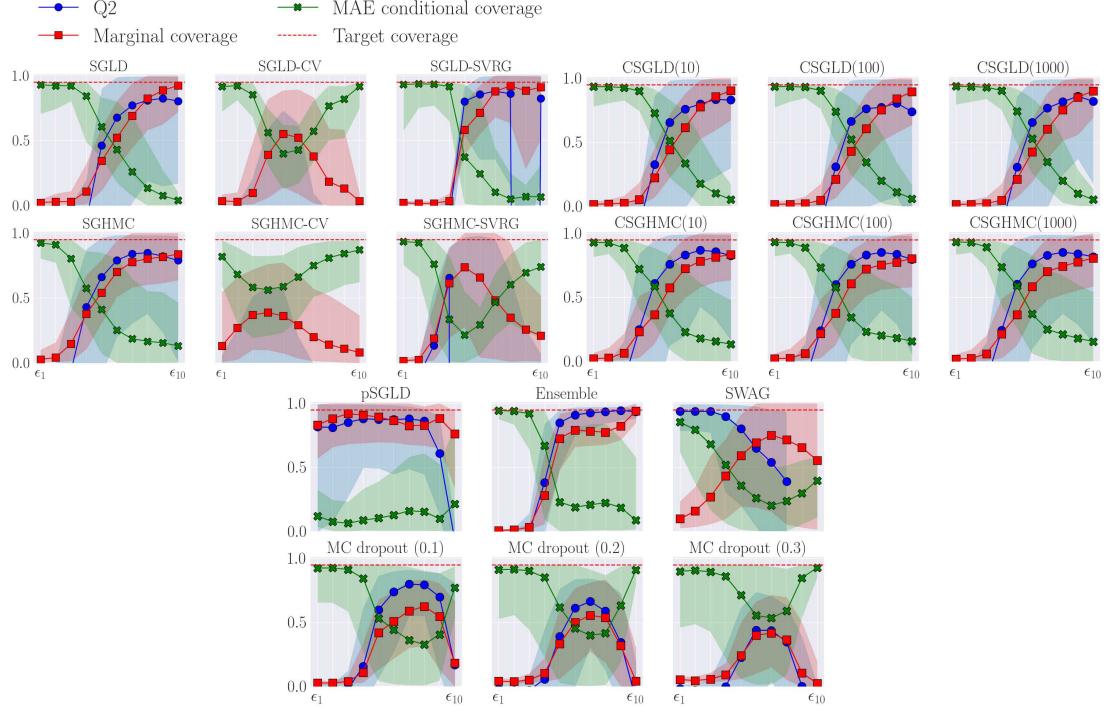


Figure 19. Problem AF#3. Coverage metrics and  $Q^2$  coefficient with respect to the step size  $\epsilon$ . The target coverage is set to 0.95.

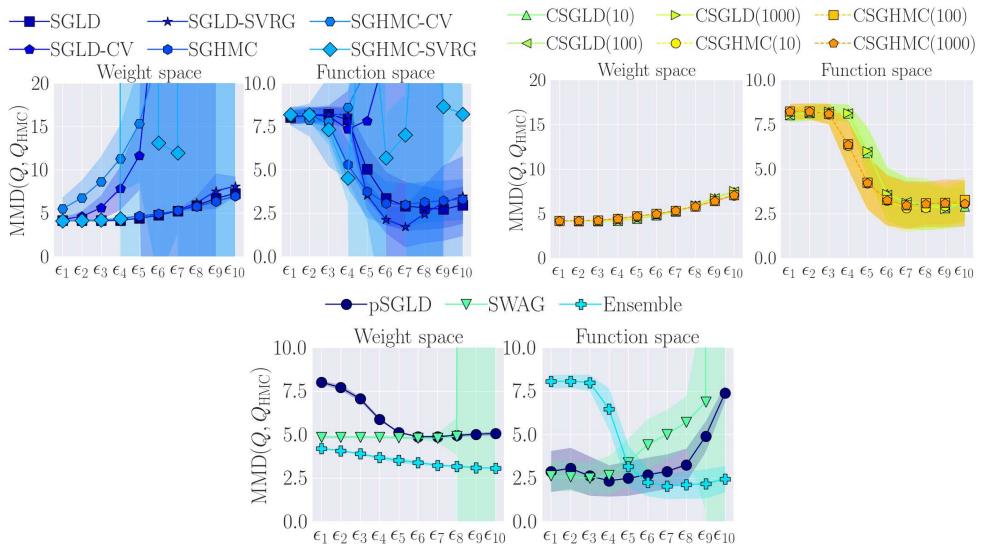


Figure 20. Problem AF#3. Maximum mean discrepancies  $MMD(Q, Q_{HMC})$  between the approximated posterior distributions  $Q$  and the reference HMC sample, in weight and function spaces.

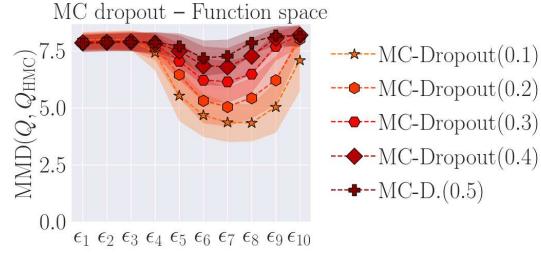


Figure 21. Problem AF#3. Maximum mean discrepancies between the approximated posterior distributions and the reference HMC sample, in weight and function spaces.

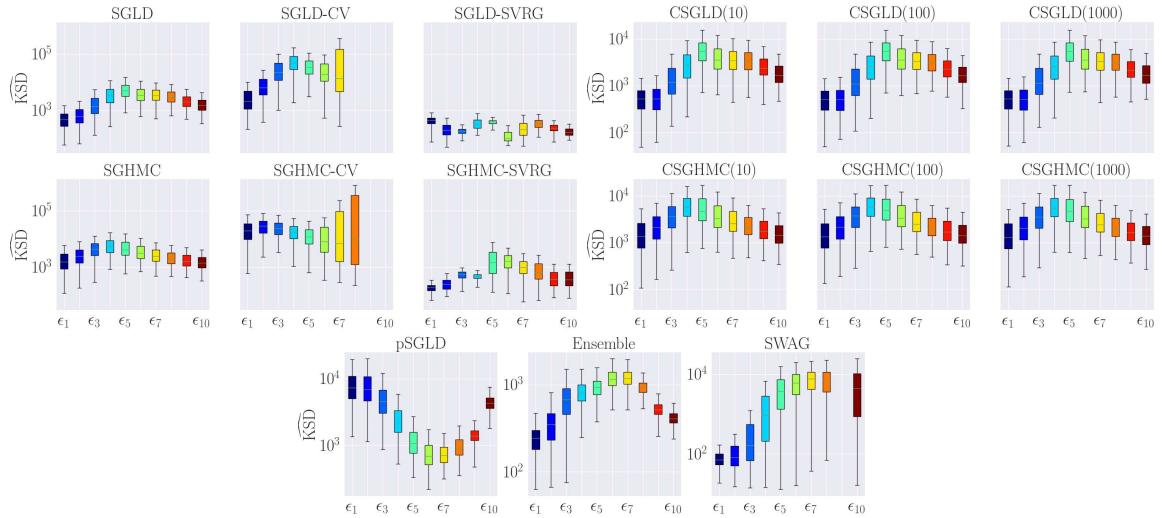


Figure 22. Problem AF#3. Kernelized Stein discrepancies  $KSD(\mathbb{P}, \mathbb{Q})$  between the target posterior measure  $\mathbb{P}$  and the approximated posteriors  $\mathbb{Q}$ .

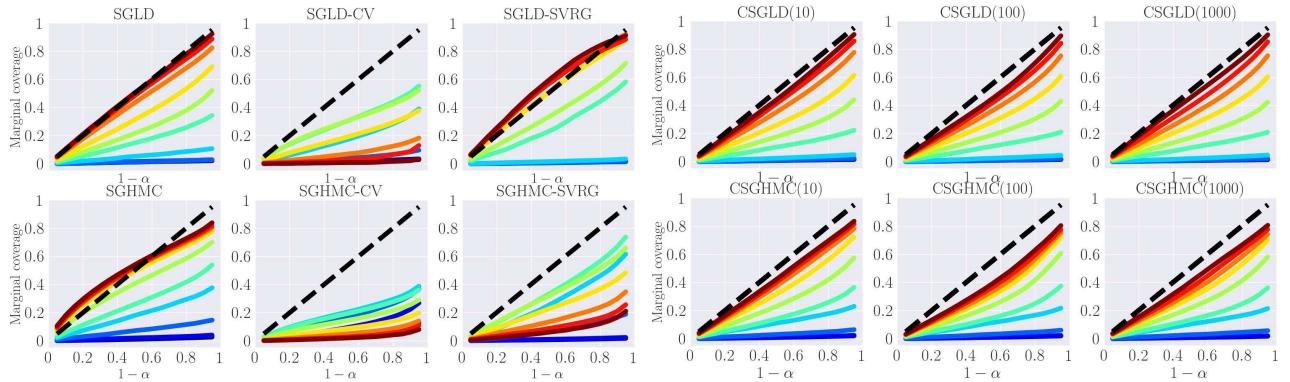


Figure 23. Problem AF#3. Graphs of the marginal coverage probability with respect to the target level for SGMCMC methods. The curves are colored by the value of the underlying step size: dark blue corresponds to the lowest, while dark red corresponds to the highest step size.

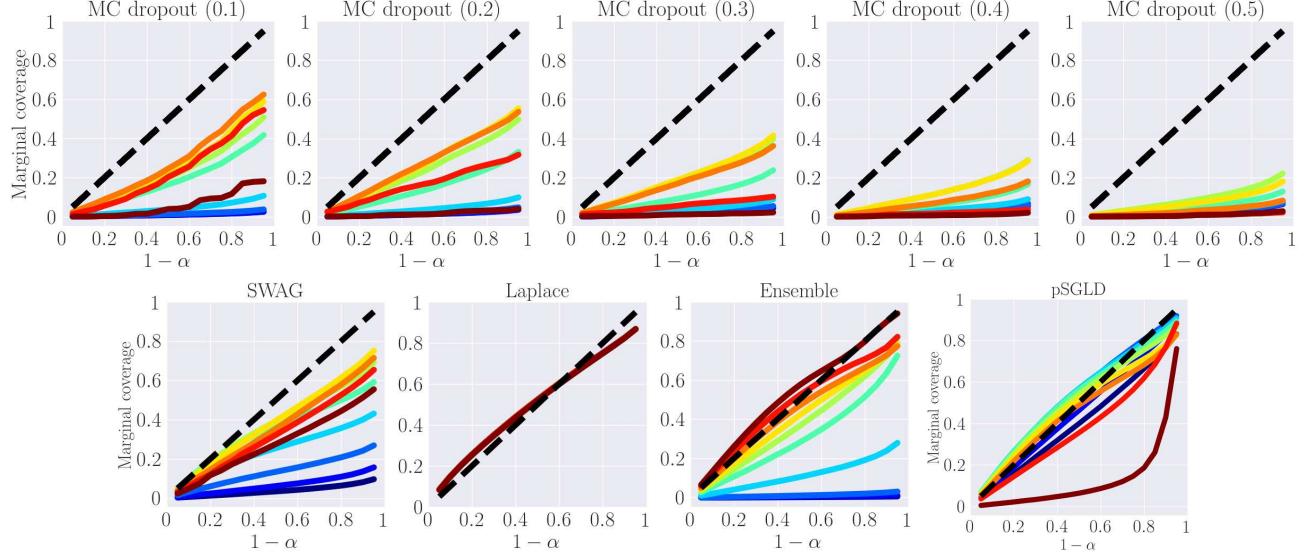


Figure 24. Problem AF#3. Graphs of the marginal coverage probability with respect to the target level for **MC-Dropout**, **SWAG**, **LA-KFAC**, **deep ensembles**, and **pSGLD**.

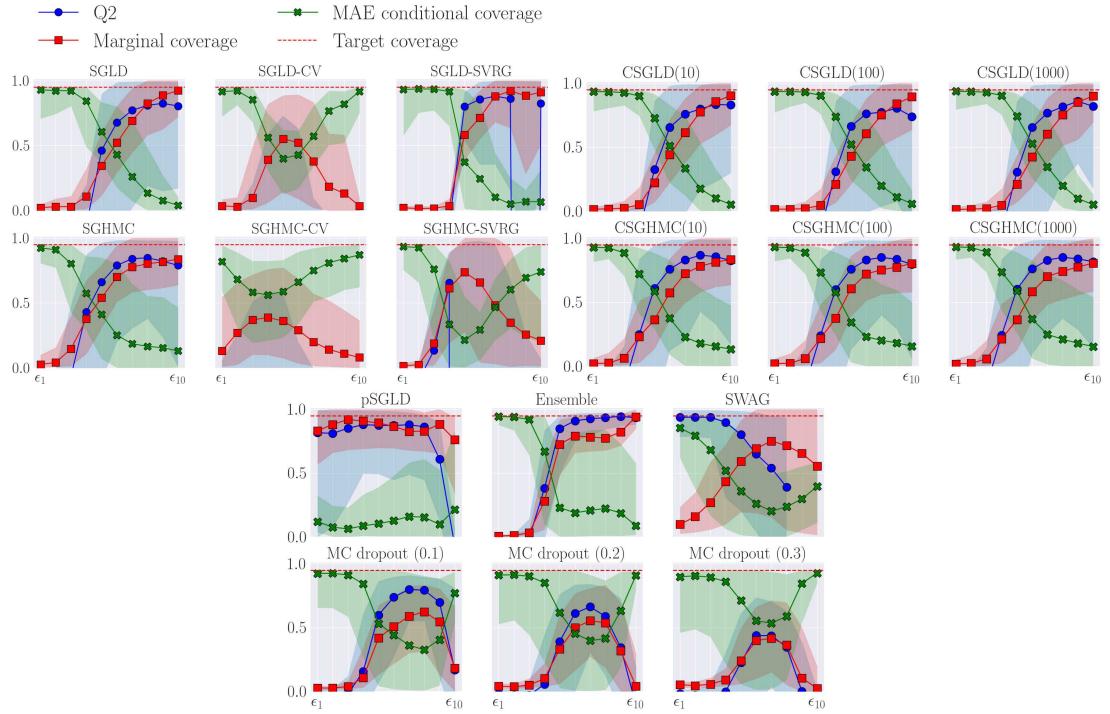


Figure 25. Problem AF#4. Coverage metrics and  $Q^2$  coefficient with respect to the step size  $\epsilon$ . The target coverage is set to 0.95.

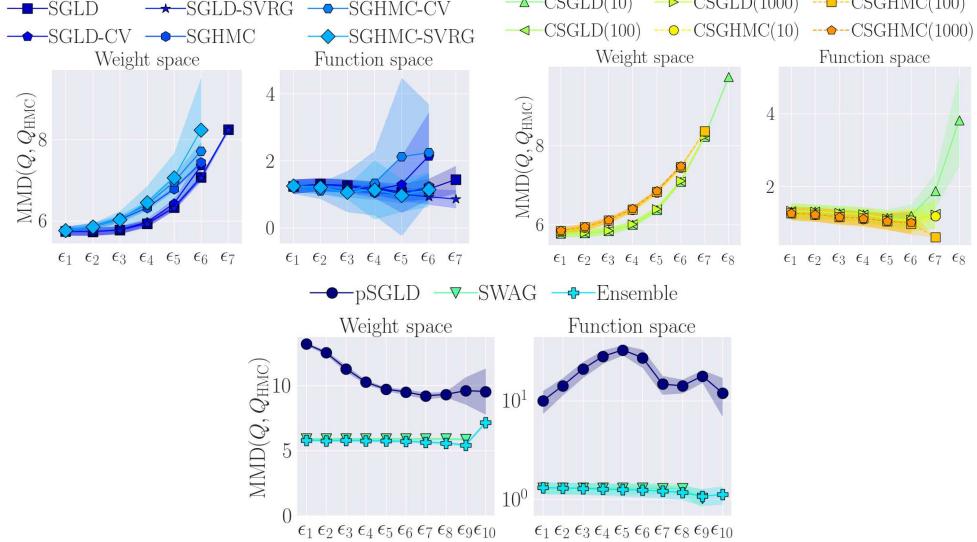


Figure 26. Problem AF#4. Maximum mean discrepancies  $\text{MMD}(\mathbb{Q}, \mathbb{Q}_{\text{HMC}})$  between the approximated posterior distributions  $\mathbb{Q}$  and the reference HMC sample, in weight and function spaces.

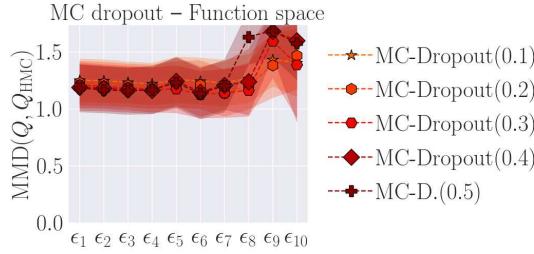


Figure 27. Problem AF#4. Maximum mean discrepancies between the approximated posterior distributions and the reference HMC sample, in weight and function spaces.

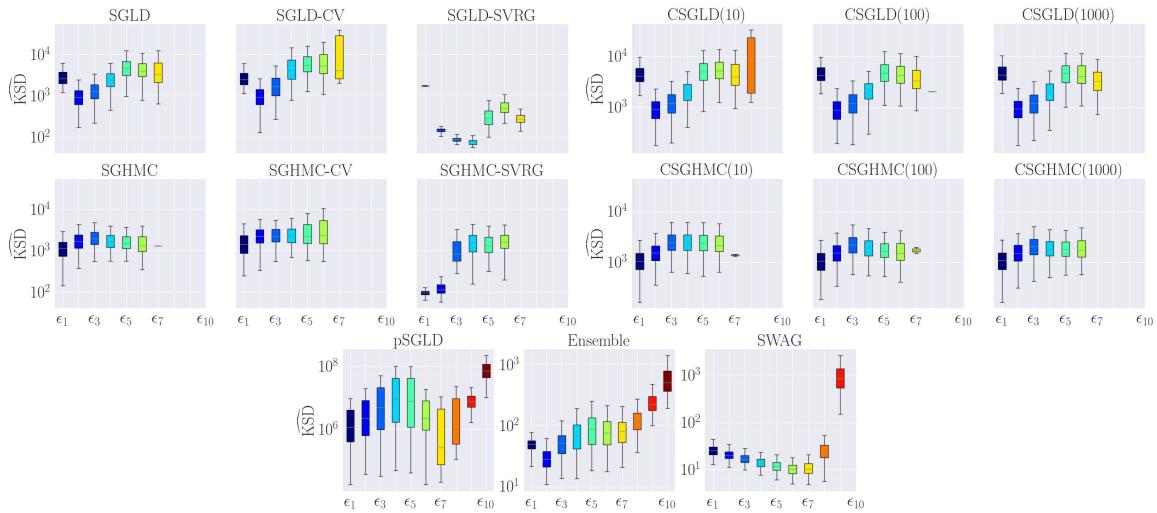


Figure 28. Problem AF#4. Kernelized Stein discrepancies  $\text{KSD}(\mathbb{P}, \mathbb{Q})$  between the target posterior measure  $\mathbb{P}$  and the approximated posteriors  $\mathbb{Q}$ .

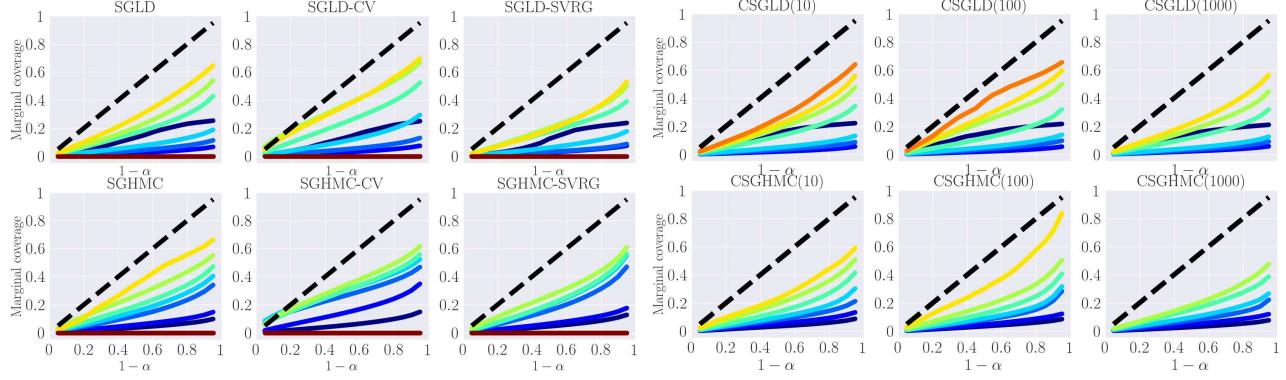


Figure 29. Problem AF#4. Graphs of the marginal coverage probability with respect to the target level for **SGMCMC** methods. The curves are colored by the value of the underlying step size: dark blue corresponds to the lowest, while dark red corresponds to the highest step size.

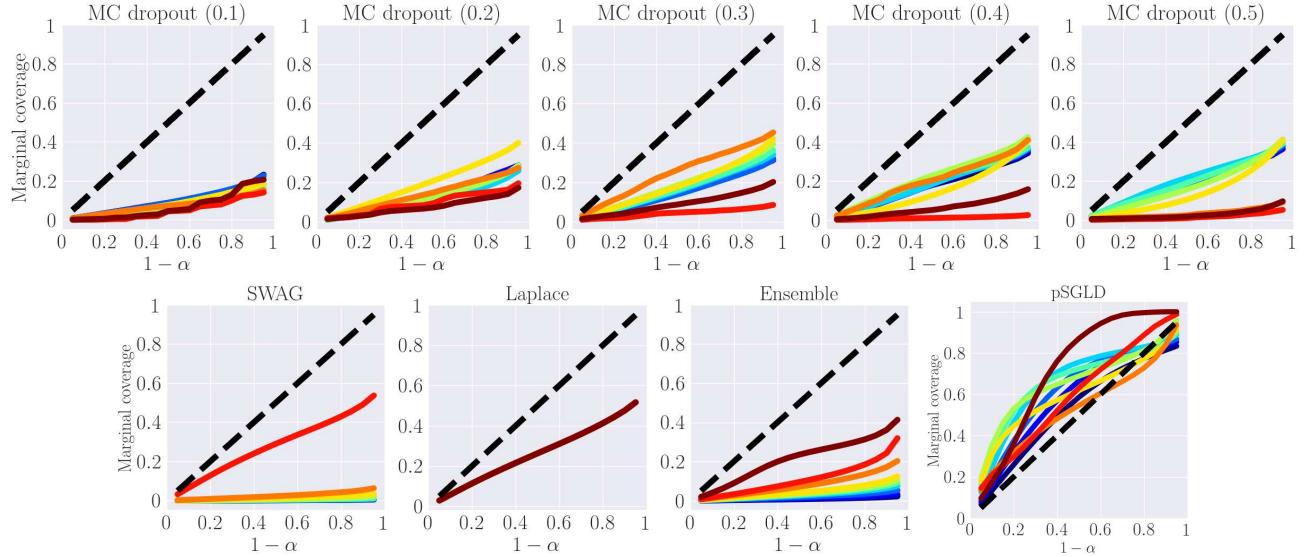


Figure 30. Problem AF#4. Graphs of the marginal coverage probability with respect to the target level for **MC-Dropout**, **SWAG**, **LA-KFAC**, **deep ensembles**, and **pSGLD**.