
Supplementary Materials for A Scalable Gradient Free Method for Bayesian Experimental Design with Implicit Models

1 Reference MI calculation using a nested Monte Carlo method

In this work, we use a nested Monte Carlo method to compute a reference mutual information (MI) value at the optimal design ξ^* . The nested Monte Carlo method provides an average of MI value and an approximation to the likelihood $p(\mathbf{y} | \xi^*, \theta)$.

Assuming D independent measurement at $\xi^* = [\xi_1^*, \dots, \xi_D^*]^\top$, the likelihood can be written by:

$$p(\mathbf{y} | \xi^*, \theta) = \prod_{j=1}^D p(y_j | \xi_j^*, \theta). \quad (1)$$

Therefore we can use Eq (1) and a sample-average of the marginal $p(\mathbf{y} | \xi^*)$ to approximate the MI:

$$I_{\text{MI}}(\xi^*) \approx \frac{1}{N_1} \sum_{i=1}^{N_1} \left[\log \frac{\prod_{j=1}^D p(y_j^{(i)} | \xi_j^*, \theta^{(i)})}{1/N_2 \sum_{k=1}^{N_2} \prod_{j=1}^D p(y_j^{(i)} | \xi_j^*, \theta^{(k)})} \right] \quad (2)$$

where $y_j^{(i)}$ are data drawn from $p(y_j | \xi_j^*, \theta^{(i)})$, $\theta^{(i)}$ and $\theta^{(k)}$ are prior samples drawn from $p(\theta)$.

1.1 Noisy linear model

For the noisy linear model, the source of noise are $\epsilon \sim \mathcal{N}(0, 1)$ and $\nu \sim \Gamma(2, 2)$, so the probability distribution of $\epsilon + \nu$ can be approximated via kernel density estimate (KDE). We generate 50,000 random samples of ϵ and ν , and compute the noise distribution $p_{\text{noise}}(\epsilon + \nu)$ using KDE with Gaussian kernels. Then we approximate the probability density $p(y_j | \xi_j^*, \theta)$ based on the sampling path for the linear model, which is given by: $\mathbf{y} = \theta_1 \mathbf{1} + \theta_2 \xi + \epsilon + \nu$. The sampling path can be reformulated to $\mathbf{y} - (\theta_1 \mathbf{1} + \theta_2 \xi) = \epsilon + \nu$. Thus we have $p(y_j | \xi_j, \theta) = p_{\text{noise}}(y_j - (\theta_0 + \theta_1 \xi))$, which enables us to compute the reference MI using Eq. (2). To obtain an accurate approximation, we use $N_1 = 10,000$ and $N_2 = 1000$ samples for the nested Monte Carlo method.

1.2 Pharmacokinetic model

For the PK model, we aim to compute the reference MI value at $\xi^* = [t_1^*, \dots, t_D^*]^\top$. Since the noise source is Gaussian noise $\epsilon_{1t} \sim \mathcal{N}(0, 0.01)$ and $\epsilon_{2t} \sim \mathcal{N}(0, 0.1)$, we can obtain an explicit equation for the data-generating distribution $p(y_j | t_j, \theta)$,

$$p(y_j | t_j, \theta) = \mathcal{N}(y_j; z(t_j, \theta), z(t_j, \theta)^2 0.01^2 + 0.1^2) \quad (3)$$

where $\theta = [V, k_a, k_e]^\top$ and the forward function $z(t_j, \theta)$ is

$$z(t_j, \theta) = \frac{D_v}{V} \frac{k_a}{k_a - k_e} [e^{-k_e t_j} - e^{-k_a t_j}]. \quad (4)$$

We can therefore approximate the reference MI value in Eq. (1) using a sample average of $p(\mathbf{y} | \xi^*)$ and data-generating distribution $p(y_j | t_j, \theta)$ in Eq. (3).

1.3 Tuning for quantum control

This example is an implicit model without pathwise gradients and no available sampling path is provided. It is infeasible to compute or approximate the reference MI value, even though the nested Monte Carlo works for the first two examples. The gradient-based BED method is not applicable for this case so we only focus on the comparison between our proposed SAGABED and two-stage method with Bayesian optimization.

2 Hyperparameter optimization

We conduct a hyperparameter optimization to search for an optimal neural network architecture and training parameters for maximizing the MI lower bound. We build a validation dataset by generating random samples from the prior distribution and data-generating distribution at the optimal design, then run five trials for every trained neural network model and estimate the mean and standard deviation. The optimal hyperparameters are obtained by means of grid search based on a set of candidate parameters.

2.1 Noisy linear model

For the low-dimensional cases $D = 1$ and $D = 10$, we consider to use one layer neural network model with different number of neurons, $\mathcal{H} = \{50, 100, 150, 200, 250, 300\}$. For the high-dimensional cases, $D = 50$ and $D = 100$, we use multiple layers from 1 to 6, and each layer has $\mathcal{H} = \{10, 30, 50, 100\}$. The learning rates is chosen properly from $\ell_\psi = \{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}\}$ and $\ell_\xi = \{10^{-3}, 10^{-2}, 10^{-1}\}$.

2.2 Pharmacokinetic model

In the PK model, we consider to use one layer neural network model with $\mathcal{H} = \{50, 100, 150, 200, 250, 300\}$ for low-dimensional case $D = 10$. For the high-dimensional cases, $D = 50$, $D = 100$ and $D = 500$. we use multiple layers from 1 to 6, and each layer has $\mathcal{H} = \{10, 30, 50, 100\}$. The candidate list of learning rates is $\ell_\psi = \{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}\}$ and $\ell_\xi = \{10^{-3}, 10^{-2}, 10^{-1}\}$.

2.3 Tuning for quantum control

We use to use one layer neural network model with $\mathcal{H} = \{50, 100, 150, 200, 250, 300\}$ for $D = 1$, $D = 5$ and $D = 10$. For the high-dimensional cases, $D = 50$, $D = 100$ and $D = 500$. we use multiple layers from 1 to 6, and each layer has $\mathcal{H} = \{10, 30, 50, 100\}$. The candidate list of learning rates is $\ell_\psi = \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$ and $\ell_\xi = \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$.

3 Additional details for the Guide ES algorithm

We use the implementation of the Guided ES (GES) algorithm published at <https://github.com/brain-research/guided-evolutionary-strategies> by the authors of the GES method. The hyperparameters of GES are listed as follows:

- α : is a hyperparameter that trades off variance between the full parameter space and the subspace. If $\alpha = 1$, it is the vanilla ES estimator, but we choose $\alpha < 1$ that can result in significantly improved performance. We use $\alpha = 0.5$ for all examples.
- β : is a scale parameter that is an overall scale factor to estimate the gradient. We use $\beta = 2.0$ for all examples.
- σ : is the variance of the perturbations or called a smooth radius. We use $\sigma = 0.1$ for the linear model and PK examples, and $\sigma = 0.5$ for quantum control example.
- P : us the number of sample pairs or called population size. It depends on the dimensionality of the problems. So we use $P = 100$ for low-dimensional cases, such as $D = 1, 10$, and we use $P = 500$ instead for high-dimensional cases, for example, $D = 50, 100, 500$.