

A Computable Measure of Suboptimality for Entropy-Regularised Variational Objectives

Clémentine Chazal, Heishiro Kanagawa, Zheyang Shen,
Anna Korba, Chris. J. Oates

October 22th, 2025

Introduction

Consider a $P \in \mathcal{P}(\mathbb{R}^d)$ as the minimizer of

$$\mathcal{J}(Q) := \mathcal{L}(Q) + \text{KLD}(Q||Q_0) \quad (1)$$

where Q_0 is a reference distribution and \mathcal{L} a loss on $\mathcal{P}(\mathbb{R}^d)$.

Introduction

Consider a $P \in \mathcal{P}(\mathbb{R}^d)$ as the minimizer of

$$\mathcal{J}(Q) := \mathcal{L}(Q) + \text{KLD}(Q||Q_0) \quad (1)$$

where Q_0 is a reference distribution and \mathcal{L} a loss on $\mathcal{P}(\mathbb{R}^d)$.

The principal issue is that P is not tractable,

Introduction

Consider a $P \in \mathcal{P}(\mathbb{R}^d)$ as the minimizer of

$$\mathcal{J}(Q) := \mathcal{L}(Q) + \text{KLD}(Q||Q_0) \quad (1)$$

where Q_0 is a reference distribution and \mathcal{L} a loss on $\mathcal{P}(\mathbb{R}^d)$.

The principal issue is that P is not tractable,

- ▶ We do not have access to the unnormalized density of P (except if \mathcal{L} is linear: if $\mathcal{L} = \int v(x)dQ(x)$, then $\mathcal{J}(Q) = \text{KLD}(Q||e^{-v}Q_0)$ and $P \propto e^{-v}Q_0$).

Introduction

Consider a $P \in \mathcal{P}(\mathbb{R}^d)$ as the minimizer of

$$\mathcal{J}(Q) := \mathcal{L}(Q) + \text{KLD}(Q||Q_0) \quad (1)$$

where Q_0 is a reference distribution and \mathcal{L} a loss on $\mathcal{P}(\mathbb{R}^d)$.

The principal issue is that P is not tractable,

- ▶ We do not have access to the unnormalized density of P (except if \mathcal{L} is linear: if $\mathcal{L} = \int v(x)dQ(x)$, then $\mathcal{J}(Q) = \text{KLD}(Q||e^{-v}Q_0)$ and $P \propto e^{-v}Q_0$).
- ▶ \mathcal{J} cannot be computed on discrete distributions as $\hat{Q} = \sum_{i=1}^n \delta_{x_i}$ because of the Kullback Leibler term.

Introduction

Consider a $P \in \mathcal{P}(\mathbb{R}^d)$ as the minimizer of

$$\mathcal{J}(Q) := \mathcal{L}(Q) + \text{KLD}(Q \| Q_0) \quad (1)$$

where Q_0 is a reference distribution and \mathcal{L} a loss on $\mathcal{P}(\mathbb{R}^d)$.

The principal issue is that P is not tractable,

- ▶ We do not have access to the unnormalized density of P (except if \mathcal{L} is linear: if $\mathcal{L} = \int v(x) dQ(x)$, then $\mathcal{J}(Q) = \text{KLD}(Q \| e^{-v} Q_0)$ and $P \propto e^{-v} Q_0$).
- ▶ \mathcal{J} cannot be computed on discrete distributions as $\hat{Q} = \sum_{i=1}^n \delta_{x_i}$ because of the Kullback Leibler term.

Intuition : Instead of minimizing \mathcal{J} , minimizing the 'size of the gradient of \mathcal{J} ',
$$\|\nabla_v \mathcal{J}(Q)\| = \sup_{\|v\| \leq 1} \langle \nabla_v \mathcal{J}(Q), v \rangle_{L^2(Q)}.$$

Kernel Gradient Discrepancy

Gradient Discrepancy

If Q and Q_0 admit density function respectively q and q_0 ,

$$\nabla_V \mathcal{J}(Q)(x) = \nabla_V \mathcal{L}(Q)(x) + \nabla \log \frac{q(x)}{q_0(x)}.$$

Gradient Discrepancy

If Q and Q_0 admit density function respectively q and q_0 ,

$$\nabla_V \mathcal{J}(Q)(x) = \nabla_V \mathcal{L}(Q)(x) + \nabla \log \frac{q(x)}{q_0(x)}.$$

Projecting the $\nabla_V \mathcal{J}(Q)$ on the vector field $v : \mathbb{R}^d \rightarrow \mathbb{R}^d$ gives

$$\begin{aligned} & \int \nabla_V \mathcal{J}(Q)(x) \cdot v(x) \, dQ(x) \\ &= \int [\nabla_V \mathcal{L}(Q)(x) - (\nabla \log q_0)(x)] \cdot v(x) \, dQ(x) + \int (\nabla \log q)(x) \cdot v(x) \, dQ(x) \\ &= \int [\nabla_V \mathcal{L}(Q)(x) - (\nabla \log q_0)(x)] \cdot v(x) \, dQ(x) - \int (\nabla \cdot v)(x) \, dQ(x) \\ &= - \int \mathcal{T}_Q v(x) \, dQ(x), \quad \mathcal{T}_Q v(x) := [(\nabla \log q_0)(x) - \nabla_V \mathcal{L}(Q)(x)] \cdot v(x) + (\nabla \cdot v)(x). \end{aligned}$$

Gradient Discrepancy

If Q and Q_0 admit density function respectively q and q_0 ,

$$\nabla_V \mathcal{J}(Q)(x) = \nabla_V \mathcal{L}(Q)(x) + \nabla \log \frac{q(x)}{q_0(x)}.$$

Projecting the $\nabla_V \mathcal{J}(Q)$ on the vector field $v : \mathbb{R}^d \rightarrow \mathbb{R}^d$ gives

$$\begin{aligned} & \int \nabla_V \mathcal{J}(Q)(x) \cdot v(x) \, dQ(x) \\ &= \int [\nabla_V \mathcal{L}(Q)(x) - (\nabla \log q_0)(x)] \cdot v(x) \, dQ(x) + \int (\nabla \log q)(x) \cdot v(x) \, dQ(x) \\ &= \int [\nabla_V \mathcal{L}(Q)(x) - (\nabla \log q_0)(x)] \cdot v(x) \, dQ(x) - \int (\nabla \cdot v)(x) \, dQ(x) \\ &= - \int \mathcal{T}_Q v(x) \, dQ(x), \quad \mathcal{T}_Q v(x) := [(\nabla \log q_0)(x) - \nabla_V \mathcal{L}(Q)(x)] \cdot v(x) + (\nabla \cdot v)(x). \end{aligned}$$

Then, one can define the **Gradient Discrepancy**,

$$\text{GD}(Q) := \sup_{\substack{v \in \mathcal{V} \text{ s.t.} \\ (\mathcal{T}_Q v)_{-} \in \mathcal{L}^1(Q)}} \left| \int \mathcal{T}_Q v(x) \, dQ(x) \right|.$$

Kernel Gradient Discrepancy (KGD)

Let $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ be a matrix-valued kernel. Let $\mathcal{B}_K = \{v \in \mathcal{H}_K : \|v\|_{\mathcal{H}_K} \leq 1\}$. The **Kernel Gradient Discrepancy (KGD)** is defined as

$$\text{KGD}_K(Q) := \sup_{\substack{v \in \mathcal{B}_K \text{ s.t.} \\ (\mathcal{T}_Q v)_- \in \mathcal{L}^1(Q)}} \left| \int \mathcal{T}_Q v(x) \, dQ(x) \right|$$

Kernel Gradient Discrepancy (KGD)

Let $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ be a matrix-valued kernel. Let $\mathcal{B}_K = \{v \in \mathcal{H}_K : \|v\|_{\mathcal{H}_K} \leq 1\}$. The **Kernel Gradient Discrepancy (KGD)** is defined as

$$\text{KGD}_K(Q) := \sup_{\substack{v \in \mathcal{B}_K \text{ s.t.} \\ (\mathcal{T}_Q v)_- \in \mathcal{L}^1(Q)}} \left| \int \mathcal{T}_Q v(x) \, dQ(x) \right|$$

Remarking that $\text{KGD}_K(Q) = \sup_{\substack{v \in \mathcal{B}_K \text{ s.t.} \\ (\mathcal{T}_Q v)_- \in \mathcal{L}^1(Q)}} \left\langle \int k_K^Q(x, \cdot) dQ(x), v \right\rangle$ where

$$k_K^Q(x, x') := \sum_{i=1}^d \sum_{j=1}^d \frac{1}{\rho_Q(x) \rho_Q(x')} \partial_{x'_j} \partial_{x_i} (\rho_Q(x) K_{i,j}(x, x') \rho_Q(x'))$$

and $\rho_Q(x) := q_0(x) \exp(-\mathcal{L}'(Q)(x))$,

Kernel Gradient Discrepancy (KGD)

Let $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ be a matrix-valued kernel. Let $\mathcal{B}_K = \{v \in \mathcal{H}_K : \|v\|_{\mathcal{H}_K} \leq 1\}$. The **Kernel Gradient Discrepancy (KGD)** is defined as

$$\text{KGD}_K(Q) := \sup_{\substack{v \in \mathcal{B}_K \text{ s.t.} \\ (\mathcal{T}_Q v)_- \in \mathcal{L}^1(Q)}} \left| \int \mathcal{T}_Q v(x) \, dQ(x) \right|$$

Remarking that $\text{KGD}_K(Q) = \sup_{\substack{v \in \mathcal{B}_K \text{ s.t.} \\ (\mathcal{T}_Q v)_- \in \mathcal{L}^1(Q)}} \left\langle \int k_K^Q(x, \cdot) dQ(x), v \right\rangle$ where

$$k_K^Q(x, x') := \sum_{i=1}^d \sum_{j=1}^d \frac{1}{\rho_Q(x) \rho_Q(x')} \partial_{x'_j} \partial_{x_i} (\rho_Q(x) K_{i,j}(x, x') \rho_Q(x'))$$

and $\rho_Q(x) := q_0(x) \exp(-\mathcal{L}'(Q)(x))$, we finally get

$$\text{KGD}_K(Q) = \left(\iint k_K^Q(x, x') \, dQ(x) dQ(x') \right)^{1/2}.$$

Experiments

Experiments : MFNN

How to sample from P ? A popular algorithm: **MFLD** (Mean Field Langevin Dynamics algorithm),

$$X_i^{t+1} = X_i^t + \epsilon[(\nabla \log q_0) - \nabla_V \mathcal{L}(Q_n^t)](X_i^t) + \sqrt{2\epsilon}Z_t^i, \quad Z_t^i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1), \quad Q_n^t := \frac{1}{n} \sum_{j=1}^n \delta_{X_j^t},$$

Experiments : MFNN

How to sample from P ? A popular algorithm: **MFLD** (Mean Field Langevin Dynamics algorithm),

$$X_i^{t+1} = X_i^t + \epsilon[(\nabla \log q_0) - \nabla_V \mathcal{L}(Q_n^t)](X_i^t) + \sqrt{2\epsilon} Z_t^i, \quad Z_t^i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1), \quad Q_n^t := \frac{1}{n} \sum_{j=1}^n \delta_{X_j^t},$$

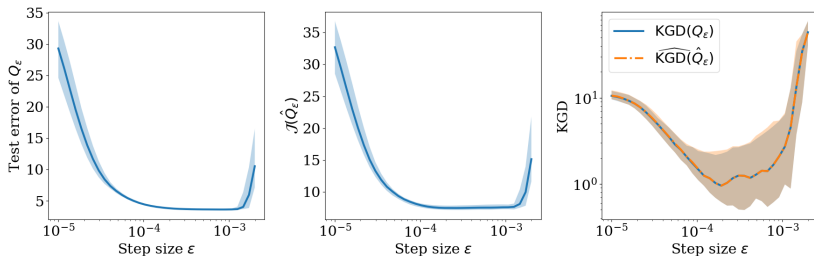
MFNN : Consider independant observations $(z_1, y_n), \dots, (z_N, y_N)$ linked by $y_i = f(z_i) + \xi_i$, $\xi_i \sim \mathcal{N}(0, \sigma^2)$ where f is a target function. We take \mathcal{L} to be the loss of a regression problem

$$\mathcal{L}(Q) = \frac{\lambda}{N} \sum_{i=1}^N \ell(y_i, \mathbb{E}_{X \sim Q}[\Phi(z_i, X)]), \quad (2)$$

where Φ is a Neural Network with parameter X . We want $f \approx \mathbb{E}_{X \sim Q}[\Phi(z_i, X)]$.

MFNN : Stepsize selection with KGD

We propose KGD as a measure to evaluate the best step size for MFLD.



MFNN : Novel sampling algorithms

For this example, we have implemented two new methods whose purpose is to optimise KGD:

- **Variational Inference:** Consider $Q_\theta = T_{\#}^\theta \mu_0$ for a reference distribution μ_0 , we solve

$$\theta_\star \in \arg \min_{\theta \in \Theta} \text{KGD}_K(Q_\theta)$$

by doing a gradient descent on θ .

MFNN : Novel sampling algorithms

For this example, we have implemented two new methods whose purpose is to optimise KGD:

- ▶ **Variational Inference:** Consider $Q_\theta = T_{\#}^\theta \mu_0$ for a reference distribution μ_0 , we solve

$$\theta_\star \in \arg \min_{\theta \in \Theta} \text{KGD}_K(Q_\theta)$$

by doing a gradient descent on θ .

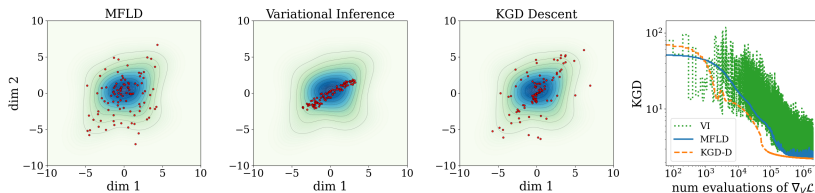
- ▶ **KGD Descent:** Let's take the discrete distribution $\hat{Q}_n = \frac{1}{n} \sum_{j=1}^n \delta_{x_j}$, we solve

$$\{x_1, \dots, x_n\} \in \arg \min \text{KGD}_K(\hat{Q}_n)$$

with gradient descent : $x_i^{t+1} = x_i^t - \varepsilon \nabla_V \text{KGD}_K^2(Q_n^t)(x_i^t)$.

MFNN : Comparison of all the methods

Here is plot the final distribution of the parameters for all the methods.



Predictively Oriented Posteriors

Another example of application is **Predictively Oriented Posteriors**. Let $p(\cdot|x)$ a parametric statistical model for independant data $\{y_i\}_{i=1,\dots,N}$. Let's take

$$\mathcal{L}(Q) = \frac{1}{2\lambda_N} \text{MMD}^2 \left(\int p(\cdot|x) \, dQ(x), \frac{1}{N} \sum_{i=1}^N \delta_{y_i} \right) \quad (3)$$

Predictively Oriented Posteriors

Another example of application is **Predictively Oriented Posteriors**. Let $p(\cdot|x)$ a parametric statistical model for independent data $\{y_i\}_{i=1,\dots,N}$. Let's take

$$\mathcal{L}(Q) = \frac{1}{2\lambda_N} \text{MMD}^2 \left(\int p(\cdot|x) \, dQ(x), \frac{1}{N} \sum_{i=1}^N \delta_{y_i} \right) \quad (3)$$

Algorithms used for this example:

- **Extensible Sampling:** Start from $x_0 \in \mathbb{R}^d$ and then apply the iterative algorithm:

$$x_n \in \arg \min_{x \in \mathbb{R}^d} \text{KGD}_K \left(\frac{1}{n} \delta_x + \frac{1}{n} \sum_{i=1}^{n-1} \delta_{x_i} \right) \quad (4)$$

where the minimum is searched on a grid in \mathbb{R}^d .

Predictively Oriented Posteriors: Variational Gradient Descent

- **Variational Gradient Descent:** This algorithm is a generalised version of SVGD.
Let $v : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $\varepsilon > 0$

$$\left. \frac{d}{d\varepsilon} \mathcal{J}((I_d + \varepsilon v)_{\#} Q) \right|_{\varepsilon=0} = - \int \mathcal{T}_Q v(x) \, dQ(x).$$

Predictively Oriented Posteriors: Variational Gradient Descent

- **Variational Gradient Descent:** This algorithm is a generalised version of SVGD. Let $v : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $\epsilon > 0$

$$\left. \frac{d}{d\epsilon} \mathcal{J}((I_d + \epsilon v)_{\#} Q) \right|_{\epsilon=0} = - \int \mathcal{T}_Q v(x) dQ(x).$$

Then the optimal direction is proportional to $\int k_K^Q(x, \cdot) dQ(x)$ which is

$$v_Q(\cdot) \propto \int \{k(x, \cdot)(\nabla \log q_0 - \nabla_V \mathcal{L}(Q))(x) + \nabla_1 k(x, \cdot)\} dQ(x),$$

Predictively Oriented Posteriors: Variational Gradient Descent

- **Variational Gradient Descent:** This algorithm is a generalised version of SVGD.
Let $v : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $\epsilon > 0$

$$\left. \frac{d}{d\epsilon} \mathcal{J}((I_d + \epsilon v)_\# Q) \right|_{\epsilon=0} = - \int \mathcal{T}_Q v(x) dQ(x).$$

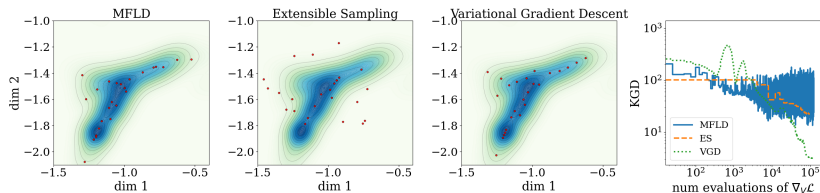
Then the optimal direction is proportional to $\int k_K^Q(x, \cdot) dQ(x)$ which is

$$v_Q(\cdot) \propto \int \{k(x, \cdot)(\nabla \log q_0 - \nabla_V \mathcal{L}(Q))(x) + \nabla_1 k(x, \cdot)\} dQ(x),$$

And then, we deduce the sampling algorithm sampling algorithm :

$$x_i^{t+1} = x_i^t + \frac{1}{n} \sum_{j=1}^n k(x_i^t, x_j^t)(\nabla \log q_0 - \nabla_V \mathcal{L}(Q_n^t))(x_j^t) + \nabla_1 k(x_j^t, x_i^t),$$

Predictively Oriented Posteriors : Comparison of the methods



Thank you !