

Optimal Thinning of MCMC Output

Chris. J. Oates
Newcastle University
Alan Turing Institute

April 2021 @ UQSay Seminar Series



Computation for the Bayesian Framework

The goal is to obtain an approximation to the posterior in a Bayesian context:

$$P : \pi(\theta|y) = \frac{\pi(y|\theta)\pi(\theta)}{\pi(y)}$$

where $\theta \in \Theta$ are the unknown parameters of the model, $\pi(\theta)$ is an appropriate prior density and y denotes the dataset.

This raises technical challenges as the normalisation constant

$$\pi(y) = \int_{\Theta} \pi(y|\theta)\pi(\theta)d\theta$$

is an intractable d -dimensional integral.

Sampling from P via Markov chain Monte Carlo (MCMC) is a popular approach which requires only evaluation of the un-normalised form

$$p(\theta) := \pi(y|\theta)\pi(\theta),$$

but it is not a silver bullet.

Computation for the Bayesian Framework

The goal is to obtain an approximation to the posterior in a Bayesian context:

$$P : \pi(\theta|y) = \frac{\pi(y|\theta)\pi(\theta)}{\pi(y)}$$

where $\theta \in \Theta$ are the unknown parameters of the model, $\pi(\theta)$ is an appropriate prior density and y denotes the dataset.

This raises technical challenges as the normalisation constant

$$\pi(y) = \int_{\Theta} \pi(y|\theta)\pi(\theta)d\theta$$

is an intractable d -dimensional integral.

Sampling from P via Markov chain Monte Carlo (MCMC) is a popular approach which requires only evaluation of the un-normalised form

$$p(\theta) := \pi(y|\theta)\pi(\theta),$$

but it is not a silver bullet.

Computation for the Bayesian Framework

The goal is to obtain an approximation to the posterior in a Bayesian context:

$$P : \pi(\theta|y) = \frac{\pi(y|\theta)\pi(\theta)}{\pi(y)}$$

where $\theta \in \Theta$ are the unknown parameters of the model, $\pi(\theta)$ is an appropriate prior density and y denotes the dataset.

This raises technical challenges as the normalisation constant

$$\pi(y) = \int_{\Theta} \pi(y|\theta)\pi(\theta)d\theta$$

is an intractable d -dimensional integral.

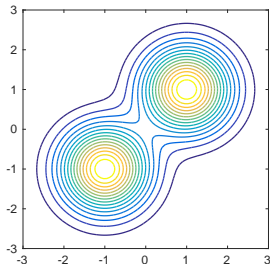
Sampling from P via Markov chain Monte Carlo (MCMC) is a popular approach which requires only evaluation of the un-normalised form

$$p(\theta) := \pi(y|\theta)\pi(\theta),$$

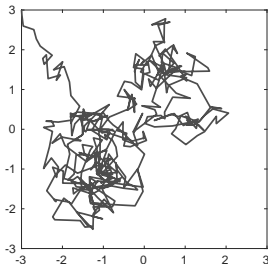
but it is not a silver bullet.

An Ideal Post-Processing Method

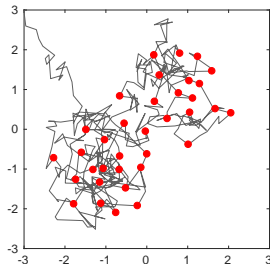
In an ideal world we would be able to post-process the MCMC output and keep only those states that are representative of the posterior P :



P



MCMC output
 $(\theta_i)_{i=1}^n$



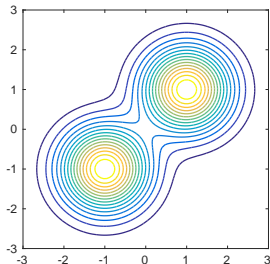
Representative Subset
 $(\theta_i)_{i \in S}$

Desiderata:

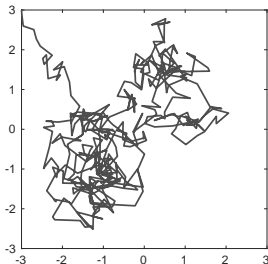
- ▶ Fix problems with MCMC (*automatic identification of burn-in; mitigation of poor mixing; number of points proportional to the probability mass in a region; etc.*)
- ▶ Compressed representation of the posterior, to reduce any downstream computational load.

An Ideal Post-Processing Method

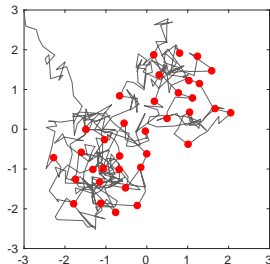
In an ideal world we would be able to post-process the MCMC output and keep only those states that are representative of the posterior P :



P



MCMC output
 $(\theta_i)_{i=1}^n$



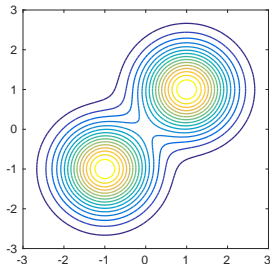
Representative Subset
 $(\theta_i)_{i \in S}$

Desiderata:

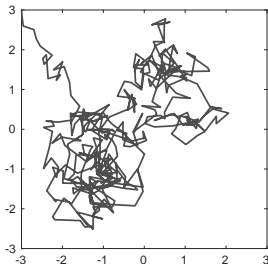
- ▶ Fix problems with MCMC (*automatic identification of burn-in; mitigation of poor mixing; number of points proportional to the probability mass in a region; etc.*)
- ▶ Compressed representation of the posterior, to reduce any downstream computational load.

An Ideal Post-Processing Method

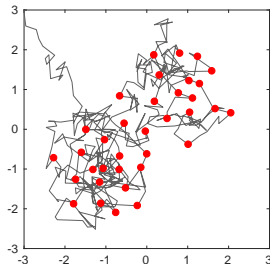
In an ideal world we would be able to post-process the MCMC output and keep only those states that are representative of the posterior P :



P



MCMC output
 $(\theta_i)_{i=1}^n$



Representative Subset
 $(\theta_i)_{i \in S}$

Desiderata:

- ▶ Fix problems with MCMC (*automatic identification of burn-in; mitigation of poor mixing; number of points proportional to the probability mass in a region; etc.*)
- ▶ Compressed representation of the posterior, to reduce any downstream computational load.

Optimal Thinning of MCMC Output

“Pick a representative subset from the MCMC output”

Idea:
$$\arg \min_{\substack{S \subset \{1, \dots, n\} \\ |S|=m}} \underbrace{\text{diff}}_{(*)} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right)$$

Remarks:

- ▶ “Nice idea, but we don’t have access to P .”
- ▶ “Combinatorial optimisation is a hard problem.”

Our strategy is to use **Stein’s Method** to manufacture a function $(*)$ that can be computed without the normalisation constant $\pi(y)$.

Optimal Thinning of MCMC Output

“Pick a representative subset from the MCMC output”

Idea:
$$\arg \min_{\substack{S \subset \{1, \dots, n\} \\ |S|=m}} \underbrace{\text{diff}}_{(*)} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right)$$

Remarks:

- ▶ “Nice idea, but we don’t have access to P .”
- ▶ “Combinatorial optimisation is a hard problem.”

Our strategy is to use **Stein’s Method** to manufacture a function $(*)$ that can be computed without the normalisation constant $\pi(y)$.

Optimal Thinning of MCMC Output

“Pick a representative subset from the MCMC output”

Idea:
$$\arg \min_{\substack{S \subset \{1, \dots, n\} \\ |S|=m}} \underbrace{\text{diff}}_{(*)} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right)$$

Remarks:

- ▶ “Nice idea, but we don’t have access to P .”
- ▶ “Combinatorial optimisation is a hard problem.”

Our strategy is to use **Stein’s Method** to manufacture a function $(*)$ that can be computed without the normalisation constant $\pi(y)$.

Optimal Thinning of MCMC Output

“Pick a representative subset from the MCMC output”

Idea:

$$\arg \min_{\substack{S \subset \{1, \dots, n\} \\ |S|=m}} \underbrace{\text{diff}}_{(*)} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right)$$

Remarks:

- ▶ “Nice idea, but we don’t have access to P .”
- ▶ “Combinatorial optimisation is a hard problem.”

Our strategy is to use **Stein’s Method** to manufacture a function $(*)$ that can be computed without the normalisation constant $\pi(y)$.

Outline

Kernel Stein Discrepancy

Stein Thinning of MCMC Output

Stein's Method in Computational Statistics

Kernel Stein Discrepancy

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition:** $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right) &:= \sup_{\|f\|_{\mathcal{K}} \leq 1} \left| \frac{1}{m} \sum_{i \in S} f(\theta_i) - \mathbb{E}_{\vartheta \sim P}[f(\vartheta)] \right| \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S})^2 = \left\| \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}}^2$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition:** $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in \mathcal{S}} \delta(\theta_i), P \right) &:= \sup_{\|f\|_{\mathcal{K}} \leq 1} \left| \frac{1}{m} \sum_{i \in \mathcal{S}} f(\theta_i) - \mathbb{E}_{\vartheta \sim P}[f(\vartheta)] \right| \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}})^2 = \left\| \frac{1}{m} \sum_{i \in \mathcal{S}} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}}^2$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition:** $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right) &:= \sup_{\|f\|_{\mathcal{K}} \leq 1} \left| \frac{1}{m} \sum_{i \in S} \langle f, k(\theta_i, \cdot) \rangle_{\mathcal{K}} - \mathbb{E}_{\vartheta \sim P} [\langle f, k(\vartheta, \cdot) \rangle_{\mathcal{K}}] \right| \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S})^2 = \left\| \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}}^2$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P,P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P,P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition: $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)**

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right) &:= \sup_{\|f\|_{\mathcal{K}} \leq 1} \left| \left\langle f, \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \mathbb{E}_{\vartheta \sim P}[k(\vartheta, \cdot)] \right\rangle_{\mathcal{K}} \right| \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S})^2 = \left\| \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}}^2$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition:** $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in \mathcal{S}} \delta(\theta_i), P \right) &:= \left\| \frac{1}{m} \sum_{i \in \mathcal{S}} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}} \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}})^2 = \left\| \frac{1}{m} \sum_{i \in \mathcal{S}} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}}^2$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition:** $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right) &:= \left\| \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}} \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S})^2 = \left\| \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}}^2$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition:** $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right) &:= \left\| \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}} \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S})^2 = \left\| \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}}^2$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition:** $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right) &:= \left\| \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}} \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in S})^2 = \left\langle \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta), \frac{1}{m} \sum_{i \in S} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\rangle_{\mathcal{K}}$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition:** $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in \mathcal{S}} \delta(\theta_i), P \right) &:= \left\| \frac{1}{m} \sum_{i \in \mathcal{S}} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}} \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$\begin{aligned} D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}})^2 &= \frac{1}{m^2} \sum_{i, j \in \mathcal{S}} \langle k(\theta_i, \cdot), k(\theta_j, \cdot) \rangle_{\mathcal{K}} - \frac{2}{m} \sum_{i \in \mathcal{S}} \int \langle k(\theta, \cdot), k(\theta_i, \cdot) \rangle_{\mathcal{K}} dP(\theta) \\ &\quad - \int \int \langle k(\theta, \cdot), k(\theta', \cdot) \rangle_{\mathcal{K}} dP(\theta) dP(\theta') \end{aligned}$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition:** $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in \mathcal{S}} \delta(\theta_i), P \right) &:= \left\| \frac{1}{m} \sum_{i \in \mathcal{S}} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}} \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}})^2 = \frac{1}{m^2} \sum_{i, j \in \mathcal{S}} k(\theta_i, \theta_j) - \frac{2}{m} \sum_{i \in \mathcal{S}} k_P(\theta_i) + k_{P, P}$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

Approximation in Reproducing Kernel Hilbert Spaces (RKHS)

Let $k : \Theta \times \Theta \rightarrow \mathbb{R}$ be the reproducing kernel of a RKHS \mathcal{K} of functions from Θ to \mathbb{R} ; i.e. $\forall \theta \in \Theta$, $k(\theta, \cdot) \in \mathcal{K}$ and $f(\theta) = \langle f, k(\theta, \cdot) \rangle_{\mathcal{K}}$ whenever $f \in \mathcal{K}$. **(Intuition: $f(\theta) = \sum_i c_i k(\theta, \theta_i)$)**

Consider an **integral probability metric** based on $\|\cdot\|_{\mathcal{K}}$:

$$\begin{aligned} \text{diff} \left(\frac{1}{m} \sum_{i \in \mathcal{S}} \delta(\theta_i), P \right) &:= \left\| \frac{1}{m} \sum_{i \in \mathcal{S}} k(\theta_i, \cdot) - \int k(\theta, \cdot) dP(\theta) \right\|_{\mathcal{K}} \\ &=: D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}}) \end{aligned}$$

which is known as the *worst-case integration error* for the RKHS \mathcal{K} .

Let's try to compute this:

$$D_{\mathcal{K}, P}(\{\theta_i\}_{i \in \mathcal{S}})^2 = \frac{1}{m^2} \sum_{i, j \in \mathcal{S}} k(\theta_i, \theta_j) - \frac{2}{m} \sum_{i \in \mathcal{S}} k_P(\theta_i) + k_{P, P}$$

where $k_P := \int k(\theta, \cdot) dP(\theta) \in \mathcal{K}$ and $k_{P, P} := \int k_P dP$.

Problem: We need to choose k carefully, so that k_P and $k_{P, P}$ can be evaluated. How?

A BOUND FOR THE ERROR IN THE
NORMAL APPROXIMATION TO THE
DISTRIBUTION OF A SUM OF
DEPENDENT RANDOM VARIABLES

CHARLES STEIN
STANFORD UNIVERSITY



Stein Characterisation

Definition (Stein Characterisation)

A distribution P is characterised by the pair $(\mathcal{A}, \mathcal{F})$, consisting of a Stein Operator \mathcal{A} and a Stein Class \mathcal{F} , if it holds that

$$\vartheta \sim P \quad \text{iff} \quad \mathbb{E}[\mathcal{A}f(\vartheta)] = 0 \quad \forall f \in \mathcal{F}.$$

Example (Stein, 1972)

- ▶ $P = N(\mu, \sigma^2)$ with density function $p(x)$
- ▶ $\mathcal{A} : f \mapsto \frac{\nabla(fp)}{p}$
- ▶ $\mathcal{F} = \{f : \mathbb{R} \rightarrow \mathbb{R} \text{ s.t. } \nabla(fp) \in L^1(\mathbb{R}) \text{ and } \lim_{x \searrow -\infty} f(\theta)p(\theta) = \lim_{\theta \nearrow +\infty} f(\theta)p(\theta)\}$.

Stein Characterisation

Definition (Stein Characterisation)

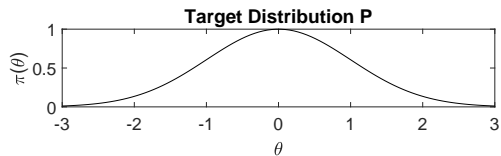
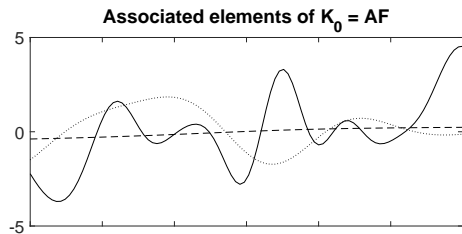
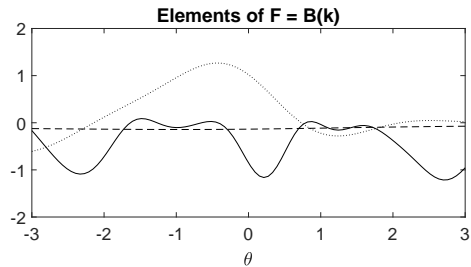
A distribution P is characterised by the pair $(\mathcal{A}, \mathcal{F})$, consisting of a Stein Operator \mathcal{A} and a Stein Class \mathcal{F} , if it holds that

$$\vartheta \sim P \quad \text{iff} \quad \mathbb{E}[\mathcal{A}f(\vartheta)] = 0 \quad \forall f \in \mathcal{F}.$$

Example (Stein, 1972)

- ▶ $P = N(\mu, \sigma^2)$ with density function $p(x)$
- ▶ $\mathcal{A} : f \mapsto \frac{\nabla(fp)}{p}$
- ▶ $\mathcal{F} = \{f : \mathbb{R} \rightarrow \mathbb{R} \text{ s.t. } \nabla(fp) \in L^1(\mathbb{R}) \text{ and } \lim_{x \searrow -\infty} f(\theta)p(\theta) = \lim_{\theta \nearrow +\infty} f(\theta)p(\theta)\}$.

Stein Characterisation



Stein Operators in Hilbert Spaces

(Going to stick to $d = 1$.)

Theorem (Chwialkowski et al. [2016])

Suppose that k is bounded, symmetric, cc-universal and satisfies $\mathbb{E}_{\vartheta \sim P}[(\Delta k(\vartheta, \vartheta))^2] < \infty$. Then P has Stein characterisation $(\mathcal{A}, \mathcal{F})$, consisting of

$$\mathcal{A}f = \frac{\nabla(f\rho)}{\rho}, \quad \mathcal{F} = \mathcal{B}(k) := \{f \in \mathcal{K} : \|f\|_{\mathcal{K}} \leq 1\}.$$

Theorem (O, Girolami and Chopin [2017])

The functions $\mathcal{A}f$ just defined are precisely the elements of the unit ball in the RKHS $\mathcal{K}_0 := \mathcal{A}\mathcal{K}$ with kernel

$$\begin{aligned} k_0(\theta, \theta') &= \nabla_{\theta} \nabla_{\theta'} k(\theta, \theta') + \frac{\nabla_{\theta} \rho(\theta)}{\rho(\theta)} \nabla_{\theta'} k(\theta, \theta') \\ &\quad + \frac{\nabla_{\theta'} \rho(\theta')}{\rho(\theta')} \nabla_{\theta} k(\theta, \theta') + \frac{\nabla_{\theta} \rho(\theta)}{\rho(\theta)} \frac{\nabla_{\theta'} \rho(\theta')}{\rho(\theta')} k(\theta, \theta'). \end{aligned}$$

In particular, under regularity conditions, $(k_0)_P = 0$ and $(k_0)_{P,P} = 0$ are trivially computed.

Solution: Use k_0 in an integral probability metric!

Stein Operators in Hilbert Spaces

(Going to stick to $d = 1$.)

Theorem (Chwialkowski et al. [2016])

Suppose that k is bounded, symmetric, cc-universal and satisfies $\mathbb{E}_{\vartheta \sim P}[(\Delta k(\vartheta, \vartheta))^2] < \infty$. Then P has Stein characterisation $(\mathcal{A}, \mathcal{F})$, consisting of

$$\mathcal{A}f = \frac{\nabla(f\rho)}{\rho}, \quad \mathcal{F} = \mathcal{B}(k) := \{f \in \mathcal{K} : \|f\|_{\mathcal{K}} \leq 1\}.$$

Theorem (O, Girolami and Chopin [2017])

The functions $\mathcal{A}f$ just defined are precisely the elements of the unit ball in the RKHS $\mathcal{K}_0 := \mathcal{A}\mathcal{K}$ with kernel

$$\begin{aligned} k_0(\theta, \theta') &= \nabla_{\theta} \nabla_{\theta'} k(\theta, \theta') + \frac{\nabla_{\theta} \rho(\theta)}{\rho(\theta)} \nabla_{\theta'} k(\theta, \theta') \\ &\quad + \frac{\nabla_{\theta'} \rho(\theta')}{\rho(\theta')} \nabla_{\theta} k(\theta, \theta') + \frac{\nabla_{\theta} \rho(\theta)}{\rho(\theta)} \frac{\nabla_{\theta'} \rho(\theta')}{\rho(\theta')} k(\theta, \theta'). \end{aligned}$$

In particular, under regularity conditions, $(k_0)_P = 0$ and $(k_0)_{P,P} = 0$ are trivially computed.

Solution: Use k_0 in an integral probability metric!

Stein Operators in Hilbert Spaces

(Going to stick to $d = 1$.)

Theorem (Chwialkowski et al. [2016])

Suppose that k is bounded, symmetric, cc-universal and satisfies $\mathbb{E}_{\vartheta \sim P}[(\Delta k(\vartheta, \vartheta))^2] < \infty$. Then P has Stein characterisation $(\mathcal{A}, \mathcal{F})$, consisting of

$$\mathcal{A}f = \frac{\nabla(f\rho)}{\rho}, \quad \mathcal{F} = \mathcal{B}(k) := \{f \in \mathcal{K} : \|f\|_{\mathcal{K}} \leq 1\}.$$

Theorem (O, Girolami and Chopin [2017])

The functions $\mathcal{A}f$ just defined are precisely the elements of the unit ball in the RKHS $\mathcal{K}_0 := \mathcal{A}\mathcal{K}$ with kernel

$$\begin{aligned} k_0(\theta, \theta') &= \nabla_{\theta} \nabla_{\theta'} k(\theta, \theta') + \frac{\nabla_{\theta} \rho(\theta)}{\rho(\theta)} \nabla_{\theta'} k(\theta, \theta') \\ &\quad + \frac{\nabla_{\theta'} \rho(\theta')}{\rho(\theta')} \nabla_{\theta} k(\theta, \theta') + \frac{\nabla_{\theta} \rho(\theta)}{\rho(\theta)} \frac{\nabla_{\theta'} \rho(\theta')}{\rho(\theta')} k(\theta, \theta'). \end{aligned}$$

In particular, under regularity conditions, $(k_0)_P = 0$ and $(k_0)_{P,P} = 0$ are trivially computed.

Solution: Use k_0 in an integral probability metric!

Stein Operators in Hilbert Spaces

(Going to stick to $d = 1$.)

Theorem (Chwialkowski et al. [2016])

Suppose that k is bounded, symmetric, cc-universal and satisfies $\mathbb{E}_{\vartheta \sim P}[(\Delta k(\vartheta, \vartheta))^2] < \infty$. Then P has Stein characterisation $(\mathcal{A}, \mathcal{F})$, consisting of

$$\mathcal{A}f = \frac{\nabla(f\rho)}{\rho}, \quad \mathcal{F} = \mathcal{B}(k) := \{f \in \mathcal{K} : \|f\|_{\mathcal{K}} \leq 1\}.$$

Theorem (O, Girolami and Chopin [2017])

The functions $\mathcal{A}f$ just defined are precisely the elements of the unit ball in the RKHS $\mathcal{K}_0 := \mathcal{A}\mathcal{K}$ with kernel

$$\begin{aligned} k_0(\theta, \theta') &= \nabla_{\theta} \nabla_{\theta'} k(\theta, \theta') + [\nabla_{\theta} \log \rho(\theta)] \nabla_{\theta'} k(\theta, \theta') \\ &\quad + [\nabla_{\theta'} \log \rho(\theta') \nabla_{\theta}] k(\theta, \theta') + [\nabla_{\theta} \log \rho(\theta)] [\nabla_{\theta'} \log \rho(\theta')] k(\theta, \theta'). \end{aligned}$$

In particular, under regularity conditions, $(k_0)_P = 0$ and $(k_0)_{P,P} = 0$ are trivially computed.

Solution: Use k_0 in an integral probability metric!

Stein Operators in Hilbert Spaces

(Going to stick to $d = 1$.)

Theorem (Chwialkowski et al. [2016])

Suppose that k is bounded, symmetric, cc-universal and satisfies $\mathbb{E}_{\vartheta \sim P}[(\Delta k(\vartheta, \vartheta))^2] < \infty$. Then P has Stein characterisation $(\mathcal{A}, \mathcal{F})$, consisting of

$$\mathcal{A}f = \frac{\nabla(f\rho)}{\rho}, \quad \mathcal{F} = \mathcal{B}(k) := \{f \in \mathcal{K} : \|f\|_{\mathcal{K}} \leq 1\}.$$

Theorem (O, Girolami and Chopin [2017])

The functions $\mathcal{A}f$ just defined are precisely the elements of the unit ball in the RKHS $\mathcal{K}_0 := \mathcal{A}\mathcal{K}$ with kernel

$$\begin{aligned} k_0(\theta, \theta') &= \nabla_{\theta} \nabla_{\theta'} k(\theta, \theta') + [\nabla_{\theta} \log \rho(\theta)] \nabla_{\theta'} k(\theta, \theta') \\ &\quad + [\nabla_{\theta'} \log \rho(\theta') \nabla_{\theta}] k(\theta, \theta') + [\nabla_{\theta} \log \rho(\theta)] [\nabla_{\theta'} \log \rho(\theta')] k(\theta, \theta'). \end{aligned}$$

In particular, under regularity conditions, $(k_0)_P = 0$ and $(k_0)_{P,P} = 0$ are trivially computed.

Solution: Use k_0 in an integral probability metric!

Stein Operators in Hilbert Spaces

(Going to stick to $d = 1$.)

Theorem (Chwialkowski et al. [2016])

Suppose that k is bounded, symmetric, cc-universal and satisfies $\mathbb{E}_{\vartheta \sim P}[(\Delta k(\vartheta, \vartheta))^2] < \infty$. Then P has Stein characterisation $(\mathcal{A}, \mathcal{F})$, consisting of

$$\mathcal{A}f = \frac{\nabla(f\rho)}{\rho}, \quad \mathcal{F} = \mathcal{B}(k) := \{f \in \mathcal{K} : \|f\|_{\mathcal{K}} \leq 1\}.$$

Theorem (O, Girolami and Chopin [2017])

The functions $\mathcal{A}f$ just defined are precisely the elements of the unit ball in the RKHS $\mathcal{K}_0 := \mathcal{A}\mathcal{K}$ with kernel

$$\begin{aligned} k_0(\theta, \theta') &= \nabla_{\theta} \nabla_{\theta'} k(\theta, \theta') + [\nabla_{\theta} \log \rho(\theta)] \nabla_{\theta'} k(\theta, \theta') \\ &\quad + [\nabla_{\theta'} \log \rho(\theta') \nabla_{\theta}] k(\theta, \theta') + [\nabla_{\theta} \log \rho(\theta)] [\nabla_{\theta'} \log \rho(\theta')] k(\theta, \theta'). \end{aligned}$$

In particular, under regularity conditions, $(k_0)_P = 0$ and $(k_0)_{P,P} = 0$ are trivially computed.

Solution: Use k_0 in an integral probability metric!

Stein Thinning of MCMC Output

Stein Thinning of MCMC Output

“Greedy pick states θ_i from the MCMC output to minimise KSD”

The “Stein Thinning” algorithm that we propose produces a subset $S = \{i_1, \dots, i_m\} \subset \{1, \dots, n\}$ consisting of:

$$\begin{aligned}i_1 &\in \arg \max_{i \in \{1, \dots, n\}} p(\theta_i | y) \\i_m &\in \arg \min_{i \in \{1, \dots, n\}} \text{KSD} \left(\frac{1}{m} \sum_{j=1}^{m-1} \delta(\theta_{i_j}) + \frac{1}{m} \delta(\theta_i), P \right), \quad m \geq 2 \\&= \arg \min_{i \in \{1, \dots, n\}} \sum_{j=1}^{m-1} k_0(\theta_i, \theta_{i_j}) + \frac{k_0(\theta_i, \theta_i)}{2}\end{aligned}$$

This requires searching over a finite set only and can therefore be exactly implemented. The cost of selecting the m th point is $O(mn)$.

Stein Thinning of MCMC Output

“Greedy pick states θ_i from the MCMC output to minimise KSD”

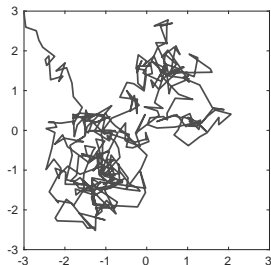
The “Stein Thinning” algorithm that we propose produces a subset $S = \{i_1, \dots, i_m\} \subset \{1, \dots, n\}$ consisting of:

$$\begin{aligned} i_1 &\in \arg \max_{i \in \{1, \dots, n\}} p(\theta_i | y) \\ i_m &\in \arg \min_{i \in \{1, \dots, n\}} \text{KSD} \left(\frac{1}{m} \sum_{j=1}^{m-1} \delta(\theta_{i_j}) + \frac{1}{m} \delta(\theta_i), P \right), \quad m \geq 2 \\ &= \arg \min_{i \in \{1, \dots, n\}} \sum_{j=1}^{m-1} k_0(\theta_i, \theta_{i_j}) + \frac{k_0(\theta_i, \theta_i)}{2} \end{aligned}$$

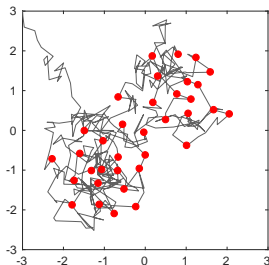
This requires searching over a finite set only and can therefore be exactly implemented. The cost of selecting the m th point is $O(mn)$.

Stein Thinning of MCMC Output

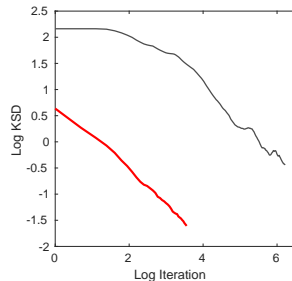
The figures we saw before were actually produced by Stein Thinning!



MCMC output
 $(\theta_i)_{i=1}^n$



Representative Subset
 $(\theta_i)_{i \in S}$



Performance
 $m \mapsto \text{KSD} \left(\frac{1}{m} \sum_{i \in S} \delta(\theta_i), P \right)$
(log-scales used)

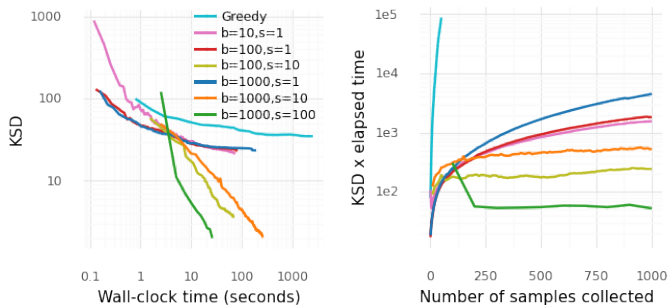
The MCMC need not even be P -invariant; full details in:

- M. Riabiz, W. Y. Chen, J. Cockayne, P. Swietach, S. A. Niederer, L. Mackey and CJO. Optimal Thinning of MCMC Output. *arXiv:2005.03952*, 2020.

Non-Myopic and Batch Extensions

However, greedy selection may be sub-optimal. Also, the cost of selecting m points from n using Stein Thinning is high, at $O(m^2n)$.

- ▶ A **non-myopic** algorithm selects s points simultaneously.
- ▶ A **mini-batch** algorithm searches over a subset of $b \ll n$ candidates at each step.



Full details in:

- ▶ O. Teymur, J. Gorham, M. Riabiz, CJO. Optimal Quantisation of Probability Measures Using Maximum Mean Discrepancy. *AISTATS*, 2021.

Stein's Method in Computational Statistics

Stein's Method in Computational Statistics

Some other uses of Stein's method in facilitating Bayesian computation:

- ▶ **Stein Points:** Chen et al. [2018, 2019]
- ▶ **Stein Importance Sampling:** Liu and Lee [2017], Hodgkinson et al. [2020]
- ▶ **Stein Variational Gradient Descent:** Liu and Wang [2016], Detommaso et al. [2018], Duncan et al. [2019], ...
- ▶ **Control Variates:** CJO et al. [2017], South et al. [2020], Si et al. [2020], ...
- ▶ **Variational Inference:** Fisher et al. [2020], ...

Recent advances in Stein discrepancies:

- ▶ **Diffusion-based Stein Operators:** Gorham and Mackey [2015], Gorham et al. [2019]
- ▶ **Stochastic Stein Discrepancy:** Huggins and Mackey [2018], Gorham et al. [2020]

Stein's Method in Computational Statistics

Some other uses of Stein's method in facilitating Bayesian computation:

- ▶ **Stein Points:** Chen et al. [2018, 2019]
- ▶ **Stein Importance Sampling:** Liu and Lee [2017], Hodgkinson et al. [2020]
- ▶ **Stein Variational Gradient Descent:** Liu and Wang [2016], Detommaso et al. [2018], Duncan et al. [2019], ...
- ▶ **Control Variates:** CJO et al. [2017], South et al. [2020], Si et al. [2020], ...
- ▶ **Variational Inference:** Fisher et al. [2020], ...

Recent advances in Stein discrepancies:

- ▶ **Diffusion-based Stein Operators:** Gorham and Mackey [2015], Gorham et al. [2019]
- ▶ **Stochastic Stein Discrepancy:** Huggins and Mackey [2018], Gorham et al. [2020]

References

- W. Chen, L. Mackey, J. Gorham, F. Briol, and CJO. Stein points. In *ICML*, 2018.
- W. Y. Chen, A. Barp, F. X. Briol, J. Gorham, L. Mackey, and CJO. Stein point Markov chain Monte Carlo. In *ICML*, 2019.
- K. Chwialkowski, H. Strathmann, and A. Gretton. A kernel test of goodness of fit. In *ICML*, 2016.
- CJO, M. Girolami, and N. Chopin. Control functionals for Monte Carlo integration. *Journal of the Royal Statistical Society: Series B*, 79(3):695–718, 2017.
- G. Detommaso, T. Cui, Y. Marzouk, R. Scheichl, and A. Spantini. A Stein variational Newton method. In *NeurIPS*, 2018.
- A. Duncan, N. Nuesken, and L. Szpruch. On the geometry of Stein variational gradient descent. *arXiv:1912.00894*, 2019.
- M. A. Fisher, T. Nolan, M. M. Graham, D. Prangle, and CJO. Measure transport with kernel Stein discrepancy. *AISTATS*, 2020.
- J. Gorham and L. Mackey. Measuring sample quality with Stein’s method. In *NeurIPS*, 2015.
- J. Gorham and L. Mackey. Measuring Sample Quality with Kernels. In *ICML*, 2017.
- J. Gorham, A. B. Duncan, S. J. Vollmer, and L. Mackey. Measuring sample quality with diffusions. *Annals of Applied Probability*, 29(5):2884–2928, 2019.
- J. Gorham, A. Raj, and L. Mackey. Stochastic Stein discrepancies. In *NeurIPS*, 2020.
- L. Hodgkinson, R. Salomone, and F. Roosta. The reproducing Stein kernel approach for post-hoc corrected sampling. *arXiv:2001.09266*, 2020.
- J. Huggins and L. Mackey. Random feature Stein discrepancies. In *NeurIPS*, 2018.
- Q. Liu and J. D. Lee. Black-box importance sampling. In *AISTATS*, 2017.
- Q. Liu and D. Wang. Stein variational gradient descent: A general purpose Bayesian inference algorithm. In *NeurIPS*, 2016.
- Q. Liu, J. Lee, and M. Jordan. A kernelized Stein discrepancy for goodness-of-fit tests. In *ICML*, 2016.
- S. Si, CJO, A. B. Duncan, L. Carin, F.-X. Briol, et al. Scalable control variates for Monte Carlo methods via stochastic optimization. *arXiv:2006.07487*, 2020.
- L. F. South, T. Karvonen, C. Nemeth, M. Girolami, and CJO. Semi-exact control functionals from Sard’s method. *arXiv:2002.00033*, 2020.