Comparison of Markov chains via weak Poincaré inequalities

Andi Q. Wang

University of Warwick

Joint with: Christophe Andrieu, Anthony Lee, Sam Power.

Joint Bayes4Health/CoSInES Workshop

March 2023









1 Introduction: Bayesian inference on modern datasets

- 2 Weak Poincaré inequalities
- 3 Application to pseudo-marginal MCMC
- 4 References

Suppose we have some (potentially vast) dataset $y = \{y_1, y_2, \dots, y_N\}$.

Suppose we have some (potentially vast) dataset $y = \{y_1, y_2, \dots, y_N\}$.

Posit a model (density function) $f_x(y)$ which generated y, which depends upon parameters $x \in \mathcal{X} = \mathbb{R}^d$.

Suppose we have some (potentially vast) dataset $y = \{y_1, y_2, \dots, y_N\}$.

Posit a model (density function) $f_x(y)$ which generated y, which depends upon parameters $x \in \mathcal{X} = \mathbb{R}^d$.

Seek learn or infer values of the parameter x which are commensurate with the observed dataset y.

Encode prior beliefs into a prior distribution $\nu(x)$, and define likelihood $\ell_y(x) := f_x(y)$.

The Bayesian approach

Encode prior beliefs into a prior distribution $\nu(x)$, and define likelihood $\ell_y(x) := f_x(y)$.

Given our observations, our posterior distribution is

$$\pi(x) = \pi(x|y) = \frac{\nu(x)\ell_y(x)}{\int \nu(z)\ell_y(z)\,\mathrm{d}z} \propto \nu(x)\ell_y(x).$$

The Bayesian approach

Encode prior beliefs into a prior distribution $\nu(x)$, and define likelihood $\ell_y(x) := f_x(y)$.

Given our observations, our posterior distribution is

$$\pi(x) = \pi(x|y) = \frac{\nu(x)\ell_y(x)}{\int \nu(z)\ell_y(z)\,\mathrm{d}z} \propto \nu(x)\ell_y(x).$$

We are then interested in quantities of the form

$$I = \pi(f) = \int_{\mathcal{X}} f(x)\pi(x) \,\mathrm{d}x,$$

e.g. $f(x) = ||x||^p$ (posterior moments), $f(x) = 1_A(x)$ (credible sets / posterior tail probabilities), etc.

Sampling

So we wish to evaluate integrals

$$I = \pi(f) = \int_{\mathcal{X}} f(x)\pi(x) \,\mathrm{d}x,$$

where π is a probability density function (our posterior distribution).

Sampling

So we wish to evaluate integrals

$$I = \pi(f) = \int_{\mathcal{X}} f(x)\pi(x) \,\mathrm{d}x,$$

where π is a probability density function (our posterior distribution).

Direct integration infeasible in high-dimensions (curse of dimensionality), furthermore only have access to π up to a normalizing constant!

Sampling

So we wish to evaluate integrals

$$I = \pi(f) = \int_{\mathcal{X}} f(x)\pi(x) \,\mathrm{d}x,$$

where π is a probability density function (our posterior distribution).

Direct integration infeasible in high-dimensions (curse of dimensionality), furthermore only have access to π up to a normalizing constant!

So instead, approximate *I* by sampling $X_1, X_2, \ldots, X_n \sim \pi$ and consider

$$I_n := \frac{1}{n} \sum_{i=1}^n f(X_i) \approx I = \int_{\mathcal{X}} f(x) \pi(x) \, \mathrm{d}x.$$

Monte Carlo

So instead, approximate I by sampling $X_1, X_2, \ldots, X_n \sim \pi$.

Monte Carlo

So instead, approximate I by sampling $X_1, X_2, \ldots, X_n \sim \pi$.

Exact sampling hard (e.g. rejection sampling also suffers from a curse of dimensionality)

So instead, approximate *I* by sampling $X_1, X_2, \ldots, X_n \sim \pi$.

Exact sampling hard (e.g. rejection sampling also suffers from a curse of dimensionality) so instead: build an ergodic Markov chain X which possesses π as its stationary distribution.

So instead, approximate *I* by sampling $X_1, X_2, \ldots, X_n \sim \pi$.

Exact sampling hard (e.g. rejection sampling also suffers from a curse of dimensionality) so instead: build an ergodic Markov chain X which possesses π as its stationary distribution.

We simulate a π -reversible ergodic Markov chain,

 X_1, X_2, \ldots

where $X_n \rightarrow \pi$ in distribution and considering

$$I_n := \frac{1}{n} \sum_{i=1}^n f(X_i) \approx I = \int_{\mathcal{X}} f(x) \pi(x) \, \mathrm{d}x.$$

Algorithm 1 (Marginal) Metropolis-Hastings (MH)

- 1: *initialise*: $X_0 = x_0, i = 0$
- 2: while i < N do
- 3: $i \leftarrow i+1$
- 4: simulate $Y_i \sim q(X_{i-1}, \cdot)$
- 5: $\alpha(X_{i-1}, Y_i) = 1 \wedge \frac{q(Y_i, X_{i-1})\pi(Y_i)}{q(X_{i-1}, Y_i)\pi(X_{i-1})}$
- 6: with probability $\alpha(X_{i-1}, Y_i)$
- 7: $X_i \leftarrow Y_i$
- 8: **else**
- 9: $X_i \leftarrow X_{i-1}$

10: return $(X_i)_{i=1,...,n}$

Algorithm 2 (Marginal) Metropolis-Hastings (MH)

- 1: *initialise*: $X_0 = x_0, i = 0$
- 2: while i < N do
- 3: $i \leftarrow i+1$
- 4: simulate $Y_i \sim q(X_{i-1}, \cdot)$
- 5: $\alpha(X_{i-1}, Y_i) = 1 \wedge \frac{q(Y_i, X_{i-1})\pi(Y_i)}{q(X_{i-1}, Y_i)\pi(X_{i-1})}$
- 6: with probability $\alpha(X_{i-1}, Y_i)$

7:
$$X_i \leftarrow Y_i$$

8: **else**

9:
$$X_i \leftarrow X_{i-1}$$

10: **return** $(X_i)_{i=1,...,n}$

We will focus on random walk Metropolis-Hastings (RWM): $q(X_{i-1}, \cdot) = \mathcal{N}(X_{i-1}, \Sigma)$.

MH example

We will be interested in studying convergence in $L^2(\pi)$:

We will be interested in studying convergence in $L^2(\pi)$:

$$L^2(\pi)=\{f:\mathcal{X} o\mathbb{R}\mid \|f\|_2^2=\int |f|^2\,\mathrm{d}\pi<\infty\},\quad \langle f,g
angle=\int fg\,\mathrm{d}\pi.$$

We will be interested in studying convergence in $L^2(\pi)$:

$$L^2(\pi) = \{f: \mathcal{X} o \mathbb{R} \mid \|f\|_2^2 = \int |f|^2 \, \mathrm{d}\pi < \infty\}, \quad \langle f, g
angle = \int fg \, \mathrm{d}\pi.$$

Given π -invariant Markov kernel P, it follows that $P: L^2(\pi) \to L^2(\pi)$ given by

 $Pf(x) = \mathbb{E}_{x}[f(X_{1})]$

is a bounded linear mapping.

We will be interested in studying convergence in $L^2(\pi)$:

$$L^2(\pi) = \{f: \mathcal{X} o \mathbb{R} \mid \|f\|_2^2 = \int |f|^2 \, \mathrm{d}\pi < \infty\}, \quad \langle f, g
angle = \int fg \, \mathrm{d}\pi.$$

Given π -invariant Markov kernel P, it follows that $P: L^2(\pi) \to L^2(\pi)$ given by

$$Pf(x) = \mathbb{E}_{x}[f(X_{1})]$$

is a bounded linear mapping.

Interested to study for $f \in L^2(\pi)$, how fast do we have

 $\|\mathbf{P}^n f - \pi(f)\|_2 \to 0.$

Exponential convergence and spectral gap

In many cases the convergence is exponential: for any $f \in L^2(\pi)$,

$$\|P^nf-\pi(f)\|_2\leq \|f\|\cdot(1-\gamma)^n,$$

where for reversible P, the best possible rate is given by the spectral gap $0 \le \gamma \le 1$.

In many cases the convergence is exponential: for any $f \in L^2(\pi)$,

$$\|P^n f - \pi(f)\|_2 \le \|f\| \cdot (1-\gamma)^n$$

where for reversible P, the best possible rate is given by the spectral gap $0 \le \gamma \le 1$.

Theorem ([Andrieu, Lee, Power, Wang (2022a)], c.f. [Roberts, Gelman, Gilks (1997)])

When the target π has a strongly concave and L-smooth potential, the spectral gap of RWM with proposal variance of d^{-1} on \mathbb{R}^d scales like

$$\gamma_d = O(d^{-1}).$$

In many cases the convergence is exponential: for any $f \in L^2(\pi)$,

```
\|P^n f - \pi(f)\|_2 \le \|f\| \cdot (1-\gamma)^n,
```

where for reversible P, the best possible rate is given by the spectral gap $0 \le \gamma \le 1$.

Theorem ([Andrieu, Lee, Power, Wang (2022a)], c.f. [Roberts, Gelman, Gilks (1997)])

When the target π has a strongly concave and L-smooth potential, the spectral gap of RWM with proposal variance of d^{-1} on \mathbb{R}^d scales like

$$\gamma_d = O(d^{-1}).$$

This is nice when applicable, but many chains actually converge at a subgeometric rate and have 0 spectral gap.

But recall:

 $\pi(x) \propto \nu(x) \ell_y(x).$

But recall:

 $\pi(x) \propto \nu(x) \ell_y(x).$

For modern datasets this $\ell_y(x)$ is often intractable! E.g.

But recall:

 $\pi(x) \propto \nu(x) \ell_y(x).$

For modern datasets this $\ell_y(x)$ is often intractable! E.g.

$$\ell_{\mathbf{y}}(\mathbf{x}) = \prod_{i=1}^{N} f_{\mathbf{x}}(y_i), \quad \text{or } \ell_{\mathbf{y}}(\mathbf{x}) = \int \int \cdots \int g(\mathbf{x}, \mathbf{z}, \mathbf{y}) \, \mathrm{d}z_1 \, \mathrm{d}z_2 \dots \mathrm{d}z_N,$$

corresponding to 'big data' or latent variable models.

$$\mathbb{E}[\hat{\ell}_y(x)] = C \cdot \ell_y(x), \quad \hat{\ell}_y(x) \ge 0$$
 a.s

$$\mathbb{E}[\hat{\ell}_y(x)] = C \cdot \ell_y(x), \quad \hat{\ell}_y(x) \ge 0$$
 a.s.

and then in Metropolis-Hastings, naively substitute in

 $\hat{\pi}(x) = \nu(x)\hat{\ell}_y(x)$

wherever we needed $\pi(x)$.

$$\mathbb{E}[\hat{\ell}_y(x)] = C \cdot \ell_y(x), \quad \hat{\ell}_y(x) \ge 0 \text{ a.s.}$$

and then in Metropolis-Hastings, naively substitute in

 $\hat{\pi}(x) = \nu(x)\hat{\ell}_y(x)$

wherever we needed $\pi(x)$.

This is the pseudo-marginal MH algorithm proposed in [Andrieu and Roberts (2009)].

Pseudo-marginal MH [Andrieu and Roberts (2009)]

Algorithm 3 Pseudo-marginal Metropolis–Hastings

- 1: *initialise*: $X_0 = x_0, i = 0$
- 2: while i < N do
- 3: $i \leftarrow i+1$
- 4: simulate $Y_i \sim q(X_{i-1}, \cdot)$
- 5: $\alpha(X_{i-1}, Y_i) = 1 \wedge \frac{q(Y_i, X_{i-1})\hat{\pi}(Y_i)}{q(X_{i-1}, Y_i)\hat{\pi}(X_{i-1})}$
- 6: with probability $\alpha(X_{i-1}, Y_i)$
- 7: $X_i \leftarrow Y_i$
- 8: **else**
- 9: $X_i \leftarrow X_{i-1}$

10: **return** $(X_i)_{i=1,...,n}$

Pseudo-marginal MH [Andrieu and Roberts (2009)]

Algorithm 4 Pseudo-marginal Metropolis-Hastings

- 1: *initialise*: $X_0 = x_0, i = 0$
- 2: while *i* < *N* do
- 3: $i \leftarrow i+1$
- 4: simulate $Y_i \sim q(X_{i-1}, \cdot)$
- 5: $\alpha(X_{i-1}, Y_i) = 1 \wedge \frac{q(Y_i, X_{i-1})\hat{\pi}(Y_i)}{q(X_{i-1}, Y_i)\hat{\pi}(X_{i-1})}$
- 6: with probability $\alpha(X_{i-1}, Y_i)$
- 7: $X_i \leftarrow Y_i$
- 8: **else**
- 9: $X_i \leftarrow X_{i-1}$

10: return $(X_i)_{i=1,...,n}$

A miracle: with an unbiased estimator, this is still a valid algorithm which targets π .

Pseudo-marginal RWM example
Natural question: how much worse?

Natural question: how much worse?

Will consider two particular examples.

• Approximate Bayesian Computation (ABC) [Marin et. al. (2012)] for Simulation-Based Inference (SBI);

Natural question: how much worse?

Will consider two particular examples.

- Approximate Bayesian Computation (ABC) [Marin et. al. (2012)] for Simulation-Based Inference (SBI);
- Particle Marginal Metropolis-Hastings (PMMH) [Andrieu et. al. (2010)] for inference with time series in State Space Models (SMMs), also known has Hidden Markov Models (HMMs).

Recall we had posterior

 $\pi(x) \propto \nu(x)\ell_y(x).$

Recall we had posterior

 $\pi(x) \propto \nu(x)\ell_y(x).$

However, $\ell_y(x) = f_x(y)$ is intractable but we can simulate from $f_x(\cdot)$.

Recall we had posterior

 $\pi(x) \propto \nu(x)\ell_y(x).$

However, $\ell_y(x) = f_x(y)$ is intractable but we can simulate from $f_x(\cdot)$.

So instead we use an approximate likelihood:

 $egin{aligned} \pi_{ ext{ABC}}(x) \propto
u(x) \ell_{ ext{ABC}}(x), \ \ell_{ ext{ABC}}(x) &\coloneqq \mathbb{P}_x(|Z-y| < \epsilon), \quad Z \sim f_x. \end{aligned}$

Recall we had posterior

$$\pi(x) \propto \nu(x)\ell_y(x).$$

However, $\ell_y(x) = f_x(y)$ is intractable but we can simulate from $f_x(\cdot)$.

So instead we use an approximate likelihood:

$$egin{aligned} \pi_{ ext{ABC}}(x) \propto
u(x) \ell_{ ext{ABC}}(x), \ \ell_{ ext{ABC}}(x) &\coloneqq \mathbb{P}_x(|Z-y| < \epsilon), \quad Z \sim f_x. \end{aligned}$$

This is further approximated using an unbiased estimator obtained by drawing M sets of artificial data $Z_1, \ldots, Z_M \stackrel{\text{i.i.d.}}{\sim} f_x$.

Recall we had posterior

 $\pi(x) \propto \nu(x)\ell_y(x).$

However, $\ell_y(x) = f_x(y)$ is intractable but we can simulate from $f_x(\cdot)$.

So instead we use an approximate likelihood:

$$egin{aligned} \pi_{ ext{ABC}}(x) \propto
u(x) \ell_{ ext{ABC}}(x), \ \ell_{ ext{ABC}}(x) &:= \mathbb{P}_x(|Z-y| < \epsilon), \quad Z \sim f_x. \end{aligned}$$

This is further approximated using an unbiased estimator obtained by drawing M sets of artificial data $Z_1, \ldots, Z_M \stackrel{\text{i.i.d.}}{\sim} f_x$.

Natural question for practitioner: what is the convergence rate and effect of varying M?

Recall we had posterior

 $\pi(x) \propto \nu(x)\ell_y(x).$

However, $\ell_y(x) = f_x(y)$ is intractable but we can simulate from $f_x(\cdot)$.

So instead we use an approximate likelihood:

$$egin{aligned} \pi_{ ext{ABC}}(x) \propto
u(x) \ell_{ ext{ABC}}(x), \ \ell_{ ext{ABC}}(x) &\coloneqq \mathbb{P}_x(|Z-y| < \epsilon), \quad Z \sim f_x. \end{aligned}$$

This is further approximated using an unbiased estimator obtained by drawing M sets of artificial data $Z_1, \ldots, Z_M \stackrel{\text{i.i.d.}}{\sim} f_x$.

Natural question for practitioner: what is the convergence rate and effect of varying M? N.B. the chain is almost always subgeometric [Lee and Łatuszyński (2014)].

In state space models, we assume there is a latent Markov chain (X_n) which drives the observation process (Y_n) . (Image: Wiki.)



In state space models, we assume there is a latent Markov chain (X_n) which drives the observation process (Y_n) . (Image: Wiki.)



The likelihood for static parameters x is typically intractable.

In state space models, we assume there is a latent Markov chain (X_n) which drives the observation process (Y_n) . (Image: Wiki.)



The likelihood for static parameters x is typically intractable.

However we can use a particle filter to get an unbiased estimator of $\ell_y(x)$; Particle Marginal Metropolis–Hastings (PMMH) [Andrieu et. al. (2010)].

In state space models, we assume there is a latent Markov chain (X_n) which drives the observation process (Y_n) . (Image: Wiki.)



The likelihood for static parameters x is typically intractable.

However we can use a particle filter to get an unbiased estimator of $\ell_y(x)$; Particle Marginal Metropolis–Hastings (PMMH) [Andrieu et. al. (2010)].

Natural questions: (subgeometric) convergence rate, how to tune the particle filter?

Andi Q. Wang (Warwick)

All pseudo-marginal examples replace $\pi(x)$ with an unbiased estimator $\hat{\pi}(x)$.

All pseudo-marginal examples replace $\pi(x)$ with an unbiased estimator $\hat{\pi}(x)$.

Since there is an extra layer of randomization, we expect pseudo-marginal MH to perform worse than the original marginal MH targeting π .

All pseudo-marginal examples replace $\pi(x)$ with an unbiased estimator $\hat{\pi}(x)$.

Since there is an extra layer of randomization, we expect pseudo-marginal MH to perform worse than the original marginal MH targeting π .

This raises the natural question: how do we quantify the degradation in performance when using pseudo-marginal MH compared to marginal MH?

All pseudo-marginal examples replace $\pi(x)$ with an unbiased estimator $\hat{\pi}(x)$.

Since there is an extra layer of randomization, we expect pseudo-marginal MH to perform worse than the original marginal MH targeting π .

This raises the natural question: how do we quantify the degradation in performance when using pseudo-marginal MH compared to marginal MH?

A first answer was given in [Andrieu and Vihola (2015)]: model $\hat{\pi}(x) = W_x \cdot \pi(x)$, with $W_x \sim Q_x$ nonnegative and $\mathbb{E}[W_x] = 1$.

All pseudo-marginal examples replace $\pi(x)$ with an unbiased estimator $\hat{\pi}(x)$.

Since there is an extra layer of randomization, we expect pseudo-marginal MH to perform worse than the original marginal MH targeting π .

This raises the natural question: how do we quantify the degradation in performance when using pseudo-marginal MH compared to marginal MH?

A first answer was given in [Andrieu and Vihola (2015)]: model $\hat{\pi}(x) = W_x \cdot \pi(x)$, with $W_x \sim Q_x$ nonnegative and $\mathbb{E}[W_x] = 1$.

Theorem ([Andrieu and Vihola (2015)])

If the marginal MH is geometric, and the W_x are uniformly bounded, then the pseudo-marginal chain is geometric. If the W_x have unbounded support, then the chain is subgeometric.

Andi Q. Wang (Warwick)

For pseudo-marginal MCMC, we are able to precisely and transparently relate the degradation in performance compared to the marginal chain with the tail probabilities of the perturbations.

For pseudo-marginal MCMC, we are able to precisely and transparently relate the degradation in performance compared to the marginal chain with the tail probabilities of the perturbations.

For ABC, we can bound the subgeometric convergence rate of the chain depending on the regularity of the prior and describe the effect of varying M.

For pseudo-marginal MCMC, we are able to precisely and transparently relate the degradation in performance compared to the marginal chain with the tail probabilities of the perturbations.

For ABC, we can bound the subgeometric convergence rate of the chain depending on the regularity of the prior and describe the effect of varying M.

For PMMH, in the limiting case of lognormal weights, we bound the subgeometric convergence rate and give precise tuning advice to minimize the asymptotic variance.

For pseudo-marginal MCMC, we are able to precisely and transparently relate the degradation in performance compared to the marginal chain with the tail probabilities of the perturbations.

For ABC, we can bound the subgeometric convergence rate of the chain depending on the regularity of the prior and describe the effect of varying M.

For PMMH, in the limiting case of lognormal weights, we bound the subgeometric convergence rate and give precise tuning advice to minimize the asymptotic variance.

Introduction: Bayesian inference on modern datasets

2 Weak Poincaré inequalities

3 Application to pseudo-marginal MCMC

4 References

Strong Poincaré inequalities

We work on
$$L^2(\mu) = \{f : \mathcal{X} \to \mathbb{R} : ||f||_2^2 < \infty\}, \quad \langle f, g \rangle := \int fg \, d\mu, \\ L_0^2(\mu) := \{f \in L^2(\mu) : \mu(f) = 0\}.$$

We work on
$$L^2(\mu) = \{f : \mathcal{X} \to \mathbb{R} : \|f\|_2^2 < \infty\}, \quad \langle f, g \rangle := \int fg \, d\mu, \\ L_0^2(\mu) := \{f \in L^2(\mu) : \mu(f) = 0\}.$$

For a μ -invariant Markov transition kernel P with $L^2(\mu)$ -adjoint P^* , define the Dirichlet form $\mathcal{E}(P^*P, f)$, for $f \in L^2_0(\mu)$:

 $\mathcal{E}(P^*P, f) := \langle (I - P^*P)f, f \rangle.$

We work on
$$L^2(\mu) = \{f : \mathcal{X} \to \mathbb{R} : \|f\|_2^2 < \infty\}, \quad \langle f, g \rangle := \int fg \, d\mu, \\ L_0^2(\mu) := \{f \in L^2(\mu) : \mu(f) = 0\}.$$

For a μ -invariant Markov transition kernel P with $L^2(\mu)$ -adjoint P^* , define the Dirichlet form $\mathcal{E}(P^*P, f)$, for $f \in L^2_0(\mu)$: $\mathcal{E}(P^*P, f) := \langle (I - P^*P)f, f \rangle.$

Strong Poincaré inequality (SPI)

A SPI holds if there exists a constant $C_{\rm P} > 0$ such that for all $f \in {\rm L}^2_0(\mu)$,

 $C_{\mathrm{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

Geometric convergence / spectral gap

 $C_{\mathbf{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

 $C_{\mathbb{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

Theorem (Geometric convergence)

Under a strong Poincaré inequality, we have for all $f \in L^2_0(\mu)$, $n \in \mathbb{N}_0$,

 $||P^n f||_2^2 \le (1 - C_{\rm P})^n ||f||_2^2.$

 $C_{\mathbf{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

Theorem (Geometric convergence)

Under a strong Poincaré inequality, we have for all $f \in L^2_0(\mu)$, $n \in \mathbb{N}_0$,

 $\|P^n f\|_2^2 \leq (1 - C_{\rm P})^n \|f\|_2^2.$

Proof. Rewriting the SPI, see $\mathcal{E}(P^*P, f)$ behaves like a discrete derivative:

 $C_{\mathbb{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

Theorem (Geometric convergence)

Under a strong Poincaré inequality, we have for all $f \in L^2_0(\mu)$, $n \in \mathbb{N}_0$,

 $\|P^n f\|_2^2 \leq (1 - C_{\rm P})^n \|f\|_2^2.$

Proof. Rewriting the SPI, see $\mathcal{E}(P^*P, f)$ behaves like a discrete derivative:

$$\begin{split} C_{\rm P} \|f\|_2^2 &\leq \mathcal{E}(P^*P, f) = \|f\|_2^2 - \langle P^*Pf, f \rangle \\ &= \|f\|_2^2 - \|Pf\|_2^2 \\ \Rightarrow \|Pf\|_2^2 &\leq (1 - C_{\rm P}) \|f\|_2^2. \end{split}$$

 $C_{\mathbf{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

Theorem (Geometric convergence)

Under a strong Poincaré inequality, we have for all $f \in L^2_0(\mu)$, $n \in \mathbb{N}_0$,

 $\|P^n f\|_2^2 \leq (1 - C_{\rm P})^n \|f\|_2^2.$

Proof. Rewriting the SPI, see $\mathcal{E}(P^*P, f)$ behaves like a discrete derivative:

$$\begin{split} C_{\rm P} \|f\|_2^2 &\leq \mathcal{E}(P^*P, f) = \|f\|_2^2 - \langle P^*Pf, f \rangle \\ &= \|f\|_2^2 - \|Pf\|_2^2 \\ \Rightarrow \|Pf\|_2^2 &\leq (1 - C_{\rm P}) \|f\|_2^2. \end{split}$$

The rest is by induction. \Box

Andi Q. Wang (Warwick)

SPI to weak Poincaré inequality

$$C_{\mathrm{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$$

$$C_{\mathrm{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$$

We now generalize this to allow for subgeometric rates of convergence:

 $C_{\mathrm{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

We now generalize this to allow for subgeometric rates of convergence:

Require $\beta : (0, \infty) \to [0, \infty)$ decreasing with $\beta(s) \downarrow 0$ as $s \to \infty$ and $\Phi : L^2(\mu) \to [0, \infty]$ given by $\Phi(f) = \|f\|_{\text{osc}}^2 = (\text{ess}_{\mu} \sup f - \text{ess}_{\mu} \inf f)^2$.
$C_{\mathrm{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

We now generalize this to allow for subgeometric rates of convergence:

Require $\beta : (0, \infty) \to [0, \infty)$ decreasing with $\beta(s) \downarrow 0$ as $s \to \infty$ and $\Phi : L^2(\mu) \to [0, \infty]$ given by $\Phi(f) = \|f\|_{\text{osc}}^2 = (\text{ess}_{\mu} \sup f - \text{ess}_{\mu} \inf f)^2$.

Weak Poincaré inequality (WPI) (*c.f.* [Röckner and Wang (2001)])

A WPI holds if:

 $C_{\mathrm{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

We now generalize this to allow for subgeometric rates of convergence:

Require $\beta : (0, \infty) \to [0, \infty)$ decreasing with $\beta(s) \downarrow 0$ as $s \to \infty$ and $\Phi : L^2(\mu) \to [0, \infty]$ given by $\Phi(f) = \|f\|_{\text{osc}}^2 = (\text{ess}_{\mu} \sup f - \text{ess}_{\mu} \inf f)^2$.

Weak Poincaré inequality (WPI) (*c.f.* [Röckner and Wang (2001)])

A WPI holds if: for some such β , Φ , $\forall s > 0$, $f \in L^2_0(\mu)$,

 $\|f\|_2^2 \leq s \mathcal{E}(P^*P, f) + \beta(s)\Phi(f).$

 $C_{\mathrm{P}}\|f\|_2^2 \leq \mathcal{E}(P^*P, f).$

We now generalize this to allow for subgeometric rates of convergence:

Require $\beta : (0, \infty) \to [0, \infty)$ decreasing with $\beta(s) \downarrow 0$ as $s \to \infty$ and $\Phi : L^2(\mu) \to [0, \infty]$ given by $\Phi(f) = \|f\|_{\text{osc}}^2 = (\text{ess}_{\mu} \sup f - \text{ess}_{\mu} \inf f)^2$.

Weak Poincaré inequality (WPI) (c.f. [Röckner and Wang (2001)])

A WPI holds if: for some such β , Φ , $\forall s > 0$, $f \in L^2_0(\mu)$,

 $\|f\|_2^2 \leq s \mathcal{E}(P^*P, f) + \beta(s)\Phi(f).$

E.g. $\beta(s) = c_0 s^{-c_1}$.

Subgeometric convergence

$$\|f\|_2^2 \leq s \, \mathcal{E}(P^*P,f) + eta(s) \Phi(f), \quad orall s > 0, f \in \mathrm{L}^2_0(\mu).$$

Subgeometric convergence

$$\|f\|_2^2 \leq s \, \mathcal{E}(P^*P, f) + \beta(s) \Phi(f), \quad \forall s > 0, f \in \mathrm{L}^2_0(\mu).$$

Define

$$egin{aligned} &\mathcal{K}(u) &\coloneqq ueta(1/u), & u \geq 0, \ &\mathcal{K}^*(v) &\coloneqq \sup_{u\geq 0}\{uv-\mathcal{K}(u)\}, & v\geq 0, \ &\mathcal{F}(x) &\coloneqq \int_x^1 rac{\mathrm{d}v}{\mathcal{K}^*(v)}, & 0 < x \leq 1. \end{aligned}$$

Subgeometric convergence

$$\|f\|_2^2 \leq s \, \mathcal{E}(P^*P, f) + eta(s) \Phi(f), \quad \forall s > 0, f \in \mathrm{L}^2_0(\mu).$$

Define

$$egin{aligned} &\mathcal{K}(u) &\coloneqq ueta(1/u), & u \geq 0, \ &\mathcal{K}^*(v) &\coloneqq \sup_{u\geq 0}\{uv-\mathcal{K}(u)\}, & v\geq 0, \ &\mathcal{F}(x) &\coloneqq \int_x^1 rac{\mathrm{d}v}{\mathcal{K}^*(v)}, & 0 < x \leq 1. \end{aligned}$$

Theorem ([Andrieu, Lee, Power, Wang (2022)])

Under a weak Poincaré inequality, we have, $\forall n \in \mathbb{N}_0$, $f \in L^2_0(\mu)$,

 $||P^nf||_2^2 \leq \Phi(f)F^{-1}(n).$

Andi Q. Wang (Warwick)

Under a weak Poincaré inequality, we have, $\forall n \in \mathbb{N}_0$, $f \in L^2_0(\mu)$,

 $||P^n f||_2^2 \leq \Phi(f) F^{-1}(n).$

If $\beta(s) = c_0 s^{-c_1}$, we can bound

 $F^{-1}(n) \leq C n^{-c_1}.$

Under a weak Poincaré inequality, we have, $\forall n \in \mathbb{N}_0$, $f \in L^2_0(\mu)$,

 $\|P^nf\|_2^2 \leq \Phi(f)F^{-1}(n).$

If $\beta(s) = c_0 s^{-c_1}$, we can bound

$$F^{-1}(n) \leq C n^{-c_1}.$$

If $\beta(s) = \eta_0 \exp(-\eta_1 s^{\eta_2})$, we can bound

$$F^{-1}(n) \leq C' \exp\left(-(Cn)^{\eta_2/(1+\eta_2)}\right).$$

Under a weak Poincaré inequality, we have, $\forall n \in \mathbb{N}_0$, $f \in L^2_0(\mu)$,

 $||P^n f||_2^2 \leq \Phi(f) F^{-1}(n).$

If $\beta(s) = c_0 s^{-c_1}$, we can bound

$$\mathsf{F}^{-1}(n) \leq C n^{-c_1}.$$

If $\beta(s) = \eta_0 \exp(-\eta_1 s^{\eta_2})$, we can bound

$$F^{-1}(n) \leq C' \exp\left(-(Cn)^{\eta_2/(1+\eta_2)}\right).$$

Intuition: the faster β decays, the faster the rate of convergence.

We have discussed WPIs of the form:

$$\|f\|_2^2 \leq s\mathcal{E}(P,f) + eta(s)\Phi(f), \quad orall s > 0, \quad orall f \in \mathrm{L}^2_0(\mu).$$

We have discussed WPIs of the form:

$$\|f\|_2^2 \leq s\mathcal{E}(P,f) + eta(s) \Phi(f), \quad orall s > 0, \quad orall f \in \mathrm{L}^2_0(\mu).$$

In order to compare Markov chains, we will consider a more general form of inequalities.

General comparison inequality

For two (reversible) Markov kernels P_1 , P_2 :

$$\mathcal{E}(P_1,f) \leq s\mathcal{E}(P_2,f) + eta(s) \Phi(f), \quad orall s > 0, \quad orall f \in \mathrm{L}^2_0(\mu).$$

We have discussed WPIs of the form:

$$\|f\|_2^2 \leq s\mathcal{E}(P,f) + eta(s) \Phi(f), \quad orall s > 0, \quad orall f \in \mathrm{L}^2_0(\mu).$$

In order to compare Markov chains, we will consider a more general form of inequalities.

General comparison inequality

For two (reversible) Markov kernels P_1 , P_2 :

 $\mathcal{E}(P_1, f) \le s \mathcal{E}(P_2, f) + \frac{\beta(s)}{\beta(s)} \Phi(f), \quad \forall s > 0, \quad \forall f \in \mathrm{L}^2_0(\mu).$ (2)

Indeed, (1) is a special case of (2) where $P_1(x, dy) = \mu(dy)$ corresponds to perfect sampling.

We have discussed WPIs of the form:

$$\|f\|_2^2 \leq s\mathcal{E}(P,f) + eta(s)\Phi(f), \quad orall s>0, \quad orall f\in \mathrm{L}^2_0(\mu).$$

In order to compare Markov chains, we will consider a more general form of inequalities.

General comparison inequality

For two (reversible) Markov kernels P_1 , P_2 :

 $\mathcal{E}(P_1, f) \le s \mathcal{E}(P_2, f) + \frac{\beta(s)}{\beta(s)} \Phi(f), \quad \forall s > 0, \quad \forall f \in \mathrm{L}^2_0(\mu).$ (2)

Indeed, (1) is a special case of (2) where $P_1(x, dy) = \mu(dy)$ corresponds to perfect sampling.

Intuition: (2) gives a bound on the convergence of P_2 relative to the convergence rate of P_1 .

Let P_1, P_2 be two μ -invariant Markov kernels on $E \times F$. Assume that for any $(x, B) \subset E \times F$,

$$P_2(x, B \setminus \{x\}) \geq \int_{B \setminus \{x\}} \epsilon(x, y) P_1(x, dy),$$

for some $\epsilon : \mathsf{E}^2 \to (0,\infty)$.

Let P_1, P_2 be two μ -invariant Markov kernels on $E \times F$. Assume that for any $(x, B) \subset E \times F$,

$$P_2(x, B \setminus \{x\}) \ge \int_{B \setminus \{x\}} \epsilon(x, y) P_1(x, dy),$$

for some $\epsilon : \mathsf{E}^2 \to (0,\infty).$

Then for any s > 0, $f \in L_0^{\infty}(\mu) \subset L_0^2(\mu)$,

 $\mathcal{E}(P_1, f) \leq \mathbf{s} \mathcal{E}(P_2, f) + \beta(\mathbf{s}) \Phi(f),$

where $\beta(s) = \frac{1}{2}\mu \otimes P_1(A(s)^{\complement} \cap \{X \neq Y\})$, $A(s) := \{(x, y) \in \mathsf{E}^2 : s \in (x, y) > 1\}$, and $\Phi(f) := \|f\|_{\mathrm{osc}}^2$.

Independence Sampler: geometric case

The Independence Sampler (IS) is one the simplest MCMC methods: given target π , at each step sample proposal $Y_i \sim q$, and accept with probability

$$lpha(\mathsf{X}_{i-1},\mathsf{Y}_i) = 1 \wedge rac{q(\mathsf{X}_{i-1})\pi(\mathsf{Y}_i)}{\pi(\mathsf{X}_{i-1})q(\mathsf{Y}_i)} = 1 \wedge rac{w(\mathsf{Y}_i)}{w(\mathsf{X}_{i-1})}.$$

Independence Sampler: geometric case

The Independence Sampler (IS) is one the simplest MCMC methods: given target π , at each step sample proposal $Y_i \sim q$, and accept with probability

$$lpha(X_{i-1},Y_i)=1\wedge rac{q(X_{i-1})\pi(Y_i)}{\pi(X_{i-1})q(Y_i)}=1\wedge rac{w(Y_i)}{w(X_{i-1})}.$$

Well-known that there is a spectral gap if and only if the weights are bounded:

$$\mathsf{w}(\mathsf{x}) \coloneqq rac{\pi(\mathsf{x})}{q(\mathsf{x})} \leq M, \quad \forall \mathsf{x} \in \mathcal{X}.$$

(E.g. if so, then you can just do rejection sampling.)

Independence Sampler: geometric case

The Independence Sampler (IS) is one the simplest MCMC methods: given target π , at each step sample proposal $Y_i \sim q$, and accept with probability

$$lpha(X_{i-1},Y_i)=1\wedge rac{q(X_{i-1})\pi(Y_i)}{\pi(X_{i-1})q(Y_i)}=1\wedge rac{w(Y_i)}{w(X_{i-1})}.$$

Well-known that there is a spectral gap if and only if the weights are bounded:

$$\mathsf{W}(x) := rac{\pi(x)}{q(x)} \leq M, \quad \forall x \in \mathcal{X}.$$

(E.g. if so, then you can just do rejection sampling.) The kernel is

$$P(x, \mathrm{d}y) = q(y) \cdot 1 \wedge rac{w(y)}{w(x)} \, \mathrm{d}y + (1 - lpha(x))\delta_x(\mathrm{d}y).$$

Independence Sampler: subgeometric case

$$P(x, \mathrm{d} y) = q(y) \cdot 1 \wedge \frac{w(y)}{w(x)} \, \mathrm{d} y + (1 - \alpha(x))\delta_x(\mathrm{d} y), \quad w(x) \coloneqq \frac{\pi(x)}{q(x)}.$$

Independence Sampler: subgeometric case

$$P(x, \mathrm{d} y) = q(y) \cdot 1 \wedge \frac{w(y)}{w(x)} \, \mathrm{d} y + (1 - \alpha(x))\delta_x(\mathrm{d} y), \quad w(x) \coloneqq \frac{\pi(x)}{q(x)}.$$

Suppose the weights are unbounded, namely

$$\sup_{x\in\mathcal{X}}w(x)=\sup_{x\in\mathcal{X}}\frac{\pi(x)}{q(x)}=\infty.$$

Independence Sampler: subgeometric case

$$P(x,\mathrm{d} y)=q(y)\cdot 1\wedge rac{w(y)}{w(x)}\,\mathrm{d} y+(1-lpha(x))\delta_x(\mathrm{d} y),\quad w(x)\coloneqq rac{\pi(x)}{q(x)}.$$

Suppose the weights are unbounded, namely

$$\sup_{x\in\mathcal{X}} w(x) = \sup_{x\in\mathcal{X}} \frac{\pi(x)}{q(x)} = \infty.$$

The problematic region is where w(x) is large, or equivalently, where $w^{-1}(x)$ is close to 0.

$$P(x, \mathrm{d} y) = q(y) \cdot 1 \wedge \frac{w(y)}{w(x)} \, \mathrm{d} y + (1 - \alpha(x))\delta_x(\mathrm{d} y), \quad w(x) := \frac{\pi(x)}{q(x)}.$$

Suppose the weights are unbounded, namely

$$\sup_{x\in\mathcal{X}}w(x)=\sup_{x\in\mathcal{X}}\frac{\pi(x)}{q(x)}=\infty.$$

The problematic region is where w(x) is large, or equivalently, where $w^{-1}(x)$ is close to 0.

The general comparison theorem establishes a WPI for the Independence Sampler: can see that $P(x, y) \ge \epsilon(x, y)\pi(y)$ for an appropriate $\epsilon(x, y)$:

$$P(x, \mathrm{d} y) = q(y) \cdot 1 \wedge \frac{w(y)}{w(x)} \, \mathrm{d} y + (1 - \alpha(x))\delta_x(\mathrm{d} y), \quad w(x) := \frac{\pi(x)}{q(x)}.$$

Suppose the weights are unbounded, namely

$$\sup_{x\in\mathcal{X}} w(x) = \sup_{x\in\mathcal{X}} \frac{\pi(x)}{q(x)} = \infty.$$

The problematic region is where w(x) is large, or equivalently, where $w^{-1}(x)$ is close to 0.

The general comparison theorem establishes a WPI for the Independence Sampler: can see that $P(x, y) \ge \epsilon(x, y)\pi(y)$ for an appropriate $\epsilon(x, y)$:

$$\epsilon(x,y) = \left(w^{-1}(x) \wedge w^{-1}(y)
ight),$$

 $\mathcal{A}(s)^{\complement} = \left\{(x,y) \in \mathsf{E} imes \mathsf{E} : \left(w^{-1}(x) \wedge w^{-1}(y)
ight) < 1/s
ight\}.$

Andi Q. Wang (Warwick)

$$P(x, \mathrm{d} y) = q(y) \cdot 1 \wedge \frac{w(y)}{w(x)} \, \mathrm{d} y + (1 - \alpha(x))\delta_x(\mathrm{d} y), \quad w(x) := \frac{\pi(x)}{q(x)}.$$

$$egin{aligned} & P(x,\mathrm{d} y) = q(y) \cdot 1 \wedge rac{w(y)}{w(x)} \, \mathrm{d} y + (1-lpha(x)) \delta_x(\mathrm{d} y), \quad w(x) \coloneqq rac{\pi(x)}{q(x)}. \ & \|f\|_2^2 \leq s \, \mathcal{E}(P,f) + eta(s) \, \|f\|_{\mathrm{osc}}^2. \end{aligned}$$

Can directly deduce a WPI for the Independence Sampler using the representations:

$$egin{aligned} \mathcal{P}(x,\mathrm{d} y) &= q(y)\cdot 1\wedge rac{w(y)}{w(x)}\,\mathrm{d} y + (1-lpha(x))\delta_x(\mathrm{d} y), \quad w(x) centcolor &= rac{\pi(x)}{q(x)}. \ &\|f\|_2^2 \leq s\,\mathcal{E}(P,f) + eta(s)\,\|f\|_{\mathrm{osc}}^2. \end{aligned}$$

Can directly deduce a WPI for the Independence Sampler using the representations:

$$\|f\|_{2}^{2} = \frac{1}{2} \int \pi(x)\pi(y)[f(y) - f(x)]^{2} dx dy,$$

$$\mathcal{E}(P, f) = \frac{1}{2} \int \pi(x)\pi(y) \left(w^{-1}(x) \wedge w^{-1}(y)\right) [f(x) - f(y)]^{2} dx dy.$$

- Introduction: Bayesian inference on modern datasets
- 2 Weak Poincaré inequalities
- 3 Application to pseudo-marginal MCMC
- 4 References

Recall our goal was to precisely characterise the degradation in convergence when using a pseudo-marginal chain with $\hat{\pi}(x) = W_x \cdot \pi(x)$ compared to the marginal MH chain with $\pi(x)$.

Recall our goal was to precisely characterise the degradation in convergence when using a pseudo-marginal chain with $\hat{\pi}(x) = W_x \cdot \pi(x)$ compared to the marginal MH chain with $\pi(x)$.

Let \tilde{P} denote the pseudo-marginal RWM chain and let P denote the marginal RWM chain. (Both constructed on an augmented space containing the weights w.)

Recall our goal was to precisely characterise the degradation in convergence when using a pseudo-marginal chain with $\hat{\pi}(x) = W_x \cdot \pi(x)$ compared to the marginal MH chain with $\pi(x)$.

Let \tilde{P} denote the pseudo-marginal RWM chain and let P denote the marginal RWM chain. (Both constructed on an augmented space containing the weights w.)

Have $W_x \sim Q_x$ nonnegative with $\mathbb{E}[W_x] = 1$, and set $\tilde{\pi}_x(\mathsf{d}w) := Q_x(\mathsf{d}w)w$.

Recall our goal was to precisely characterise the degradation in convergence when using a pseudo-marginal chain with $\hat{\pi}(x) = W_x \cdot \pi(x)$ compared to the marginal MH chain with $\pi(x)$.

Let \tilde{P} denote the pseudo-marginal RWM chain and let P denote the marginal RWM chain. (Both constructed on an augmented space containing the weights w.)

Have $W_x \sim Q_x$ nonnegative with $\mathbb{E}[W_x] = 1$, and set $\tilde{\pi}_x(\mathsf{d}w) := Q_x(\mathsf{d}w)w$.

Theorem ([Andrieu, Lee, Power, Wang (2022)])

We have for any s > 0, $f \in L^2_0(\mu)$ bounded,

 $\mathcal{E}(\mathbf{P}, f) \leq s \, \mathcal{E}(\tilde{\mathbf{P}}, f) + \frac{\beta(s)}{\|f\|_{\mathrm{osc}}^2},$

where

$$eta(s) \coloneqq \int_{\mathcal{X}} \tilde{\pi}_x(w \ge s) \pi(\mathrm{d} x), \quad s > 0.$$

Andi Q. Wang (Warwick)

Pseudo-marginal derivation

The WPI for pseudo-marginal is derived in almost the same way as for the Independence Sampler!

We work on the augmented state space $\mathcal{X} \times [0, \infty)$. The pseudo-marginal kernel \tilde{P} is given by: $\mathfrak{r}(x, y)$ is the standard MH acceptance ratio,

$$ilde{P}(x,w;\mathrm{d}y,\mathrm{d}u) = \left[1 \wedge \left\{\mathfrak{r}(x,y) rac{u}{w}
ight\}
ight] q(x,\mathrm{d}y) Q_y(\mathrm{d}u) + \delta_{x,w}(\mathrm{d}y,\mathrm{d}u) \widetilde{
ho}(x,w),$$

Pseudo-marginal derivation

The WPI for pseudo-marginal is derived in almost the same way as for the Independence Sampler!

We work on the augmented state space $\mathcal{X} \times [0, \infty)$. The pseudo-marginal kernel \tilde{P} is given by: $\mathfrak{r}(x, y)$ is the standard MH acceptance ratio,

$$ilde{P}(x,w;\mathrm{d}y,\mathrm{d}u) = \left[1 \wedge \left\{\mathfrak{r}(x,y)\frac{u}{w}\right\}\right]q(x,\mathrm{d}y)Q_y(\mathrm{d}u) + \delta_{x,w}(\mathrm{d}y,\mathrm{d}u)\widetilde{
ho}(x,w),$$

which we compare to the standard MH, on the extended state space:

 $P(x,w;\mathrm{d} y,\mathrm{d} u)\left[1\wedge\mathfrak{r}(x,y)\right]q(x,\mathrm{d} y)\,u\,Q_y(\mathrm{d} u)+\delta_{x,w}(\mathrm{d} y,\mathrm{d} u)\widetilde{\rho}(x,w)$

Pseudo-marginal derivation

The WPI for pseudo-marginal is derived in almost the same way as for the Independence Sampler!

We work on the augmented state space $\mathcal{X} \times [0, \infty)$. The pseudo-marginal kernel \tilde{P} is given by: $\mathfrak{r}(x, y)$ is the standard MH acceptance ratio,

$$ilde{P}(x,w;\mathrm{d}y,\mathrm{d}u) = \left[1 \wedge \left\{\mathfrak{r}(x,y)\frac{u}{w}\right\}\right]q(x,\mathrm{d}y)Q_y(\mathrm{d}u) + \delta_{x,w}(\mathrm{d}y,\mathrm{d}u)\widetilde{
ho}(x,w),$$

which we compare to the standard MH, on the extended state space:

$$P(x, w; dy, du) \left[1 \wedge \mathfrak{r}(x, y)\right] q(x, dy) \, \underline{u} \, Q_y(du) + \delta_{x, w}(dy, du) \tilde{\rho}(x, w)$$

The problematic region is when w gets large. We take

$$\epsilon(x,w;y,u)=w^{-1}\wedge u^{-1}.$$

Application 1: Approximate Bayesian Computation (ABC)

We run an MCMC chain targeting the ABC posterior, true likelihood $\ell_y(x) = f_x(y)$, $\pi_{ABC}(x) \propto \nu(x)\ell_{ABC}(x)$, $\ell_{ABC}(x) := \mathbb{P}_x(|Z - y| < \epsilon), \quad Z \sim f_x.$

Application 1: Approximate Bayesian Computation (ABC)

We run an MCMC chain targeting the ABC posterior, true likelihood $\ell_y(x) = f_x(y)$, $\pi_{ABC}(x) \propto \nu(x)\ell_{ABC}(x)$, $\ell_{ABC}(x) := \mathbb{P}_x(|Z - y| < \epsilon), \quad Z \sim f_x.$

This is further approximated using a pseudo-marginal approach: $W_M := \frac{1}{M} \sum_{i=1}^M W_i$,

$$W_j := 1\{ |Z_j - y| < \epsilon \} / \ell_{ABC}(x), \quad Z_j \stackrel{\mathrm{iid}}{\sim} f_x.$$
Application 1: Approximate Bayesian Computation (ABC)

We run an MCMC chain targeting the ABC posterior, true likelihood $\ell_y(x) = f_x(y)$, $\pi_{ABC}(x) \propto \nu(x)\ell_{ABC}(x)$, $\ell_{ABC}(x) := \mathbb{P}_x(|Z - y| < \epsilon), \quad Z \sim f_x.$

This is further approximated using a pseudo-marginal approach: $W_M := \frac{1}{M} \sum_{j=1}^M W_j$,

$$W_j := 1\{ \left| Z_j - y \right| < \epsilon \} / \ell_{ABC}(x), \quad Z_j \stackrel{\mathrm{iid}}{\sim} f_x.$$

Theorem ([Andrieu, Lee, Power, Wang (2022)])

Suppose that $\int_{\mathcal{X}} \nu(x) \ell_{ABC}^{-(p-1)}(x) dx < \infty$. Then there is $C_{M,p} > 0$ such that for all s > 0,

$$\beta(s) = \int_{\mathcal{X}} \pi(\mathsf{d} x) \widetilde{\pi}_x(\mathcal{W}_M \ge s) \le C_{M,p} s^{-p},$$

and as $M \to \infty$, $C_{M,p} = 1 + O(1/M)$. The convergence rate for the chain is $O(n^{-p})$.

Application 2: Particle Marginal Metropolis–Hastings (PMMH)

PMMH is a well-established algorithm to perform MCMC, e.g. for state space models [Andrieu et. al. (2010)].

Application 2: Particle Marginal Metropolis-Hastings (PMMH)

PMMH is a well-established algorithm to perform MCMC, e.g. for state space models [Andrieu et. al. (2010)].

We can study the limiting case when the weights are log-normal($-\sigma^2/2, \sigma^2$), where $\sigma^2 = \sigma_0^2/N$, N number of particles [Bérard et. al. (2014)].

Application 2: Particle Marginal Metropolis–Hastings (PMMH)

PMMH is a well-established algorithm to perform MCMC, e.g. for state space models [Andrieu et. al. (2010)].

We can study the limiting case when the weights are log-normal($-\sigma^2/2, \sigma^2$), where $\sigma^2 = \sigma_0^2/N$, N number of particles [Bérard et. al. (2014)].

Theorem ([Andrieu, Lee, Power, Wang (2022)])

Assume the marginal chain satisfies a SPI. Then we have the convergence bound, $\forall n \in \mathbb{N}$,

$$F_{\rm PMMH}^{-1}(\mathbf{n}) \leq \frac{2}{C_{\rm P}} \exp\left\{-\frac{1}{2\sigma^2} \mathsf{W}^2\left(\frac{C_{\rm P}\sigma^2}{2\exp(\sigma^2/2)} \cdot \mathbf{n}\right)\right\},$$

W is the Lambert function (inverse of $x \mapsto xe^x$).

Application to Particle Marginal Metropolis-Hastings (PMMH) II

To minimize the asymptotic variance taking into account the computational cost, (*c.f.* our results for convergence bounds and mixing times), tune algorithm so that

 $\sigma pprox$ 0.973.



This technique has been particularly fruitful for analyzing pseudo-marginal MCMC methods, significantly extending [Andrieu and Vihola (2015)],

This technique has been particularly fruitful for analyzing pseudo-marginal MCMC methods, significantly extending [Andrieu and Vihola (2015)], but the general approach could be applied much more broadly.

This technique has been particularly fruitful for analyzing pseudo-marginal MCMC methods, significantly extending [Andrieu and Vihola (2015)], but the general approach could be applied much more broadly.

See also further fundamental theory in our follow-up technical report [Andrieu, Lee, Power, Wang (2022b)].

- Introduction: Bayesian inference on modern datasets
- 2 Weak Poincaré inequalities
- 3 Application to pseudo-marginal MCMC
- 4 References

Thanks for listening! I

- Andrieu, C., Doucet, A., Holenstein, R. (2010). Particle Markov chain Monte Carlo methods. J. Roy. Statist. Soc. Ser. B: Stat. Methodol., 72(3), 269342.
- Andrieu, C., Lee, A., Power, S., Wang, A. Q. (2022). Comparison of Markov chains via weak Poincaré inequalities with application to pseudo-marginal MCMC. Ann. Statist., 50(6), 3592-3618.
- Andrieu, C., Lee, A., Power, S., Wang, A. Q. (2022). Explicit convergence bounds for Metropolis Markov chains: isoperimetry, spectral gaps and profiles. https://doi.org/10.48550/arxiv.2211.08959.
- Andrieu, C., Lee, A., Power, S., Wang, A. Q. (2022). Poincaré inequalities for Markov chains: a meeting with Cheeger, Lyapunov and Metropolis. *Technical report*. https://doi.org/10.48550/arxiv.2208.05239.
- Andrieu, C., Roberts, G. O. (2009). The pseudo-marginal approach for efficient Monte Carlo computations. *Ann. Statist.*, 37(2), 697725.
 Andrieu, C., Vihola, M. (2015). Convergence properties of pseudo-marginal Markov chain Monte Carlo algorithms. *Ann. Appl.*
 - Andrieu, C., Vihola, M. (2015). Convergence properties of pseudo-marginal Markov chain Monte Carlo algorithms. *Ann. Appl. Probab.*, 25(2), 10301077.
- Baxendale, P. H. (2005). Renewal theory and computable convergence rates for geometrically ergodic Markov chains. Ann. Appl. Probab., 15(1B), 700738.
 - Bérard, J., Del Moral, P., Doucet, A. (2014). A lognormal central limit theorem for particle approximations of normalizing constants. *Electron. J. Probab.*, 19, 128.
 - Diaconis, P., Saloff-Coste, L. (1993). Comparison Theorems for Reversible Markov Chains. Ann. Appl. Probab., 3(3), 696730.

- Lee, A., Łatuszyński, K. (2014). Variance bounding and geometric ergodicity of Markov chain Monte Carlo kernels for approximate Bayesian computation. *Biometrika*, 101(3), 655671.
- Marin, J. M., Pudlo, P., Robert, C. P., Ryder, R. J. (2012). Approximate Bayesian computational methods. *Statist. Comput.*, 22(6), 11671180.
- Roberts, G. O., Gelman, A., Gilks, W. R. (1997). Weak Convergence and Optimal Scaling of random walk Metropolis algorithms. Ann. Appl. Probab., 7(1), 110120.
- Röckner, M., Wang, F.-Y. (2001). Weak Poincare Inequalities and L² Convergence Rates of Markov Semigroups. J. Funct. Anal., 185, 564603.
- Feel free to get in touch: andi.wang@warwick.ac.uk!

Fix $f \in L^2_0(\mu)$. Have that

$$\|f\|_2^2 \leq s\mathcal{E}(P^*P, f) + \frac{\beta(s)}{\phi(f)} \Phi(f), \quad \forall s > 0$$

Fix $f \in L^2_0(\mu)$. Have that

$$\|f\|_2^2 \leq s\mathcal{E}(P^*P,f) + eta(s)\Phi(f), \quad orall s>0$$

$$\Rightarrow \frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \geq \frac{\|f\|_2^2}{s\Phi(f)} - \frac{\beta(s)}{s}.$$

Fix $f \in L^2_0(\mu)$. Have that

$$\|f\|_2^2 \leq s\mathcal{E}(P^*P,f) + eta(s)\Phi(f), \quad orall s>0$$

$$\Rightarrow \frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \geq \frac{\|f\|_2^2}{s\Phi(f)} - \frac{\beta(s)}{s}.$$

Fix $f \in L^2_0(\mu)$. Have that

$$\|f\|_2^2 \leq s\mathcal{E}(P^*P,f) + eta(s)\Phi(f), \quad orall s > 0$$

$$\Rightarrow \frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \geq \frac{\|f\|_2^2}{s\Phi(f)} - \frac{\beta(s)}{s}.$$

$$\frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \geq u \cdot \frac{\|f\|_2^2}{\Phi(f)} - \mathcal{K}(u), \quad \forall u > 0.$$

Fix $f \in L^2_0(\mu)$. Have that

$$\|f\|_2^2 \leq s\mathcal{E}(P^*P,f) + eta(s)\Phi(f), \quad orall s > 0$$

$$\Rightarrow \frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \geq \frac{\|f\|_2^2}{s\Phi(f)} - \frac{\beta(s)}{s}.$$

$$\frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \ge u \cdot \frac{\|f\|_2^2}{\Phi(f)} - \mathcal{K}(u), \quad \forall u > 0.$$
$$\frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \ge \sup_{u > 0} \left\{ u \cdot \frac{\|f\|_2^2}{\Phi(f)} - \mathcal{K}(u) \right\}$$

Fix $f \in L^2_0(\mu)$. Have that

$$\|f\|_2^2 \leq s\mathcal{E}(P^*P,f) + eta(s)\Phi(f), \quad orall s > 0$$

$$\Rightarrow \frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \geq \frac{\|f\|_2^2}{s\Phi(f)} - \frac{\beta(s)}{s}.$$

$$\frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \ge u \cdot \frac{\|f\|_2^2}{\Phi(f)} - \mathcal{K}(u), \quad \forall u > 0.$$
$$\frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \ge \sup_{u > 0} \left\{ u \cdot \frac{\|f\|_2^2}{\Phi(f)} - \mathcal{K}(u) \right\} =: \mathcal{K}^*\left(\frac{\|f\|_2^2}{\Phi(f)}\right).$$

Fix $f \in L^2_0(\mu)$. Have that

$$\|f\|_2^2 \leq s\mathcal{E}(P^*P,f) + eta(s)\Phi(f), \quad orall s>0$$

$$\Rightarrow \frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \geq \frac{\|f\|_2^2}{s\Phi(f)} - \frac{\beta(s)}{s}.$$

Set u := 1/s, $K(u) := u\beta(1/u)$.

$$\frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \ge u \cdot \frac{\|f\|_2^2}{\Phi(f)} - \mathcal{K}(u), \quad \forall u > 0.$$
$$\frac{\mathcal{E}(P^*P,f)}{\Phi(f)} \ge \sup_{u>0} \left\{ u \cdot \frac{\|f\|_2^2}{\Phi(f)} - \mathcal{K}(u) \right\} =: \mathcal{K}^*\left(\frac{\|f\|_2^2}{\Phi(f)}\right).$$

Call this final inequality optimized WPI (oWPI).

Andi Q. Wang (Warwick)

Now define

$$F(x) := \int_x^1 \frac{\mathrm{d}v}{K^*(v)}, \quad x \in (0, a], \qquad \frac{h_n}{h_n} := \frac{\|P^n f\|_2^2}{\Phi(f)}.$$

Now define

$$F(x) := \int_x^1 \frac{\mathrm{d}v}{K^*(v)}, \quad x \in (0, a], \qquad \frac{h_n}{h_n} := \frac{\|P^n f\|_2^2}{\Phi(f)}.$$

Now define

$$F(x) := \int_x^1 \frac{\mathrm{d}v}{K^*(v)}, \quad x \in (0, a], \qquad \frac{h_n}{h_n} := \frac{\|P^n f\|_2^2}{\Phi(f)}.$$

$$F(h_n) - F(h_{n-1}) = \int_{h_n}^{h_{n-1}} \frac{\mathrm{d}v}{K^*(v)}$$

Now define

$$F(x) := \int_x^1 \frac{\mathrm{d}v}{K^*(v)}, \quad x \in (0, a], \qquad \frac{h_n}{h_n} := \frac{\|P^n f\|_2^2}{\Phi(f)}.$$

$$F(h_n) - F(h_{n-1}) = \int_{h_n}^{h_{n-1}} \frac{dv}{K^*(v)} \\ \ge (h_{n-1} - h_n) / K^*(h_{n-1})$$

Now define

$$F(\mathbf{x}) := \int_{\mathbf{x}}^{1} \frac{\mathrm{d}\mathbf{v}}{K^{*}(\mathbf{v})}, \quad \mathbf{x} \in (0, \mathbf{a}], \qquad \mathbf{h}_{\mathbf{n}} := \frac{\|P^{n}f\|_{2}^{2}}{\Phi(f)}.$$

$$F(h_n) - F(h_{n-1}) = \int_{h_n}^{h_{n-1}} \frac{dv}{K^*(v)}$$

$$\geq (h_{n-1} - h_n) / K^*(h_{n-1})$$

$$= \frac{\mathcal{E}(P^*P, P^{n-1}f) / \Phi(f)}{K^*(h_{n-1})}$$

Now define

$$F(x) := \int_x^1 \frac{\mathrm{d}v}{K^*(v)}, \quad x \in (0, a], \qquad h_n := \frac{\|P^n f\|_2^2}{\Phi(f)}.$$

Want to bound convergence of $h_n \rightarrow 0$.

F

$$\begin{aligned} (h_n) - F(h_{n-1}) &= \int_{h_n}^{h_{n-1}} \frac{\mathrm{d}v}{K^*(v)} \\ &\geq (h_{n-1} - h_n) / K^*(h_{n-1}) \\ &= \frac{\mathcal{E}(P^*P, P^{n-1}f) / \Phi(f)}{K^*(h_{n-1})} \\ &\geq K^*(h_{n-1}) / K^*(h_{n-1}) = 1. \end{aligned}$$

Now define

$$F(x) := \int_x^1 \frac{\mathrm{d}v}{K^*(v)}, \quad x \in (0, a], \qquad h_n := \frac{\|P^n f\|_2^2}{\Phi(f)}.$$

$$F(h_n) - F(h_{n-1}) = \int_{h_n}^{h_{n-1}} \frac{dv}{K^*(v)}$$

$$\geq (h_{n-1} - h_n) / K^*(h_{n-1})$$

$$= \frac{\mathcal{E}(P^*P, P^{n-1}f) / \Phi(f)}{K^*(h_{n-1})}$$

$$\geq K^*(h_{n-1}) / K^*(h_{n-1}) = 1. \quad (oWPI)$$

$$\Rightarrow F(h_n) - F(h_0) \geq n.$$

Now define

$$F(x) := \int_x^1 \frac{\mathrm{d}v}{K^*(v)}, \quad x \in (0, a], \qquad h_n := \frac{\|P^n f\|_2^2}{\Phi(f)}.$$

Want to bound convergence of $h_n \rightarrow 0$.

$$F(h_n) - F(h_{n-1}) = \int_{h_n}^{h_{n-1}} \frac{dv}{K^*(v)}$$

$$\geq (h_{n-1} - h_n) / K^*(h_{n-1})$$

$$= \frac{\mathcal{E}(P^*P, P^{n-1}f) / \Phi(f)}{K^*(h_{n-1})}$$

$$\geq K^*(h_{n-1}) / K^*(h_{n-1}) = 1. \quad (oWPI)$$

$$\Rightarrow F(h_n) - F(h_0) \geq n.$$

So we invert this to obtain

$$\|P^nf\|_2^2 \leq \Phi(f)F^{-1}(n). \quad \Box$$

Andi Q. Wang (Warwick)

The definition of a WPI involved $\mathcal{E}(P^*P, f)$. In general P^*P is a complex object.

The definition of a WPI involved $\mathcal{E}(P^*P, f)$. In general P^*P is a complex object.

Recall that P is reversible if $P^* = P$, and furthermore P is positive if for all $f \in L^2(\mu)$, $\langle f, Pf \rangle \ge 0$.

The definition of a WPI involved $\mathcal{E}(P^*P, f)$. In general P^*P is a complex object.

Recall that P is reversible if $P^* = P$, and furthermore P is positive if for all $f \in L^2(\mu)$, $\langle f, Pf \rangle \ge 0$.

Lemma ([Baxendale (2005)], [Andrieu and Vihola (2015)])

The following MCMC chains are positive: Random Walk Metropolis with Gaussian proposals, the Independence Sampler, and pseudo-marginal MH built from any of these chains.

The definition of a WPI involved $\mathcal{E}(P^*P, f)$. In general P^*P is a complex object.

Recall that P is reversible if $P^* = P$, and furthermore P is positive if for all $f \in L^2(\mu)$, $\langle f, Pf \rangle \ge 0$.

Lemma ([Baxendale (2005)], [Andrieu and Vihola (2015)])

The following MCMC chains are positive: Random Walk Metropolis with Gaussian proposals, the Independence Sampler, and pseudo-marginal MH built from any of these chains.

Theorem [Andrieu, Lee, Power, Wang (2022)]

For a positive kernel P, we have for all $f \in L^2(\mu)$,

 $\mathcal{E}(P,f) \leq \mathcal{E}(P^2,f).$

Therefore a WPI for P implies a WPI for $P^*P = P^2$.