Kernel Adaptive Metropolis-Hastings

Arthur Gretton,*

*Gatsby Unit, CSML, University College London

NIPS, December 2015



Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 1 / 24

Metropolis-Hastings MCMC

- Unnormalized target $\pi(x) \propto p(x)$
- Generate Markov chain with invariant distribution p
 - Initialize $x_0 \sim p_0$
 - At iteration $t \ge 0$, propose to move to state $x' \sim q(\cdot | x_t)$
 - Accept/Reject proposals based on ratio

$$\begin{aligned} \mathbf{x}_{t+1} &= & \left\{ \begin{aligned} \mathbf{x}', & \text{w.p. min} \left\{ 1, \frac{\pi(\mathbf{x}')q(\mathbf{x_t}|\mathbf{x}')}{\pi(\mathbf{x_t})q(\mathbf{x}'|\mathbf{x_t})} \right\}, \\ \mathbf{x}_t, & \text{otherwise.} \end{aligned} \right. \end{aligned}$$

• What proposal $q(\cdot|x_t)$?

Metropolis-Hastings MCMC

- Unnormalized target $\pi(x) \propto p(x)$
- Generate Markov chain with invariant distribution p
 - Initialize $x_0 \sim p_0$
 - At iteration $t \ge 0$, propose to move to state $x' \sim q(\cdot|x_t)$
 - Accept/Reject proposals based on ratio

$$x_{t+1} = \begin{cases} x', & \text{w.p. min} \left\{ 1, \frac{\pi(x')q(x_t|x')}{\pi(x_t)q(x'|x_t)} \right\}, \\ x_t, & \text{otherwise.} \end{cases}$$

- What proposal $q(\cdot|x_t)$?
 - Too narrow or broad: \rightarrow slow convergence
 - $\bullet\,$ Does not conform to support of target \to slow convergence

(日本)

Adaptive MCMC

• Adaptive Metropolis (Haario, Saksman & Tamminen, 2001): Update proposal $q_t(\cdot|x_t) = \mathcal{N}(x_t, \nu^2 \hat{\Sigma}_t)$, using estimates of the target covariance

Adaptive MCMC

• Adaptive Metropolis (Haario, Saksman & Tamminen, 2001): Update proposal $q_t(\cdot|x_t) = \mathcal{N}(x_t, \nu^2 \hat{\Sigma}_t)$, using estimates of the target covariance

Adaptive MCMC

 Adaptive Metropolis (Haario, Saksman & Tamminen, 2001): Update proposal q_t(·|x_t) = N(x_t, ν²Σ̂_t), using estimates of the target covariance



Locally miscalibrated for *strongly non-linear targets*: directions of large variance depend on the current location

12/12/2015

3 / 24

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

Motivation: Intractable & Non-linear Targets

- Previous solutions for non-linear targets: Hamiltonian Monte Carlo (HMC) or Metropolis Adjusted Langevin Algorithms (MALA) (Roberts & Stramer, 2003; Girolami & Calderhead, 2011).
- Require target gradients and second order information

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

Motivation: Intractable & Non-linear Targets

- Previous solutions for non-linear targets: Hamiltonian Monte Carlo (HMC) or Metropolis Adjusted Langevin Algorithms (MALA) (Roberts & Stramer, 2003; Girolami & Calderhead, 2011).
- Require target gradients and second order information

Our case: not even target $\pi(\cdot)$ can be computed – Pseudo-Marginal MCMC (Beaumont, 2003; Andrieu & Roberts, 2009).

(日本)

Example: when is target not computable?

 GPC model: latent process f, labels y, (with covariate matrix X), and hyperparameters θ:

$$p(\mathbf{f}, \mathbf{y}, \theta) = p(\theta)p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f})$$

 $\mathbf{f}| heta\sim\mathcal{N}(\mathbf{0},\mathcal{K}_{ heta})$ GP with covariance $\mathcal{K}_{ heta}$

Example: when is target not computable?

 GPC model: latent process f, labels y, (with covariate matrix X), and hyperparameters θ:

$$p(\mathbf{f}, \mathbf{y}, \theta) = p(\theta)p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f})$$

 $\mathbf{f}| heta\sim\mathcal{N}(\mathbf{0},\mathcal{K}_{ heta})$ GP with covariance $\mathcal{K}_{ heta}$

• Automatic Relevance Determination (ARD) covariance:

$$(\mathcal{K}_{\theta})_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}'_j | \theta) = \exp\left(-\frac{1}{2}\sum_{s=1}^d \frac{(x_{i,s} - x'_{j,s})^2}{\exp(\theta_s)}\right)$$

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 5/24

Example: when is target not computable?

• Gaussian process classification, latent process f

 $p(\theta|\mathbf{y}) \propto p(\theta)p(\mathbf{y}|\theta) = p(\theta) \int p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f},\theta)d\mathbf{f} =: \pi(\theta)$

... but cannot integrate out **f**

Example: when is target not computable?

• Gaussian process classification, latent process f

$$p(\theta|\mathbf{y}) \propto p(\theta)p(\mathbf{y}|\theta) = p(\theta) \int p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f},\theta)d\mathbf{f} =: \pi(\theta)$$

- ... but cannot integrate out **f**
- MH ratio:

$$\alpha(\theta, \theta') = \min\left\{1, \frac{p(\theta')p(\mathbf{y}|\theta')q(\theta|\theta')}{p(\theta)p(\mathbf{y}|\theta)q(\theta'|\theta)}\right\}$$

JI SAR

(3)

Example: when is target not computable?

• Gaussian process classification, latent process f

$$p(\theta|\mathbf{y}) \propto p(\theta)p(\mathbf{y}|\theta) = p(\theta) \int p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f},\theta)d\mathbf{f} =: \pi(\theta)$$

... but cannot integrate out **f**

MH ratio:

$$\alpha(\theta, \theta') = \min\left\{1, \frac{p(\theta')p(\mathbf{y}|\theta')q(\theta|\theta')}{p(\theta)p(\mathbf{y}|\theta)q(\theta'|\theta)}\right\}$$

 Filippone & Girolami, 2013 use Pseudo-Marginal MCMC: unbiased estimate of p(y|θ) via importance sampling:

$$\hat{p}(heta|\mathbf{y}) \propto p(heta) \hat{p}(\mathbf{y}| heta) pprox p(heta) rac{1}{n_{ ext{imp}}} \sum_{i=1}^{n_{ ext{imp}}} p(\mathbf{y}|\mathbf{f}^{(i)}) rac{p(\mathbf{f}^{(i)}| heta)}{Q(\mathbf{f}^{(i)})}$$

12/12/2015 6/24

4 回 > 4 回 > 4 回 > 回 回 の Q の

Example: when is target not computable?

• Gaussian process classification, latent process **f**

$$p(\theta|\mathbf{y}) \propto p(\theta)p(\mathbf{y}|\theta) = p(\theta) \int p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f},\theta)d\mathbf{f} =: \pi(\theta)$$

... but cannot integrate out f

• Estimated MH ratio:

$$\alpha(\theta, \theta') = \min\left\{1, \frac{p(\theta')\hat{\rho}(\mathbf{y}|\theta')q(\theta|\theta')}{p(\theta)\hat{\rho}(\mathbf{y}|\theta)q(\theta'|\theta)}\right\}$$

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

Example: when is target not computable?

• Gaussian process classification, latent process **f**

$$p(\theta|\mathbf{y}) \propto p(\theta)p(\mathbf{y}|\theta) = p(\theta) \int p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f},\theta)d\mathbf{f} =: \pi(\theta)$$

- ... but cannot integrate out **f**
- Estimated MH ratio:

$$\alpha(\theta, \theta') = \min\left\{1, \frac{p(\theta')\hat{\rho}(\mathbf{y}|\theta')q(\theta|\theta')}{p(\theta)\hat{\rho}(\mathbf{y}|\theta)q(\theta'|\theta)}\right\}$$

• Replacing marginal likelihood $p(y|\theta)$ with unbiased estimate $\hat{p}(y|\theta)$ still results in correct invariant distribution [Beaumont, 2003; Andrieu & Roberts, 2009]

12/12/2015 7 / 24

A = A = A = A = A = A = A

Intractable & Non-linear Target in GPC

• Sliced posterior over hyperparameters of a Gaussian Process classifier on UCI Glass dataset obtained using Pseudo-Marginal MCMC



Adaptive sampler that learns the shape of non-linear targets without gradient information?

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 8/24

Two strategies for adaptive sampling



Kameleon (Sejdinovic et al. 2014)

- Learns covariance in RKHS.
- Locally aligns to (non-linear) target covariance, gradient free.



Kernel Adaptive Hamiltonian Monte Carlo (Strathmann et al. 2015)

• Learns *global* estimate of gradient of log target density

A B K A B K

JIN NOR

The Kameleon



D. Sejdinovic, H. Strathmann, M. Lomeli, C. Andrieu, and A. Gretton, ICML 2014

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 10 / 24

EL SQA

 \bullet Capture non-linearities using linear covariance C_z in feature space $\mathcal H$

EL SQC

11 / 24

12/12/2015

Input space \mathcal{X}

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

 \bullet Capture non-linearities using linear covariance C_z in feature space $\mathcal H$



Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

 \bullet Capture non-linearities using linear covariance C_z in feature space $\mathcal H$



 \bullet Capture non-linearities using linear covariance C_z in feature space $\mathcal H$



12/12/2015 11 / 24

I SOCA

Proposal Construction Summary

- Get a chain subsample $z = \{z_i\}_{i=1}^n$
- **2** Construct an RKHS sample $f \sim \mathcal{N}(\phi(x_t), \nu^2 C_z)$
- Solution Propose x^* such that $\phi(x^*)$ is close to f (with an additional exploration term $\xi \sim \mathcal{N}(0, \gamma^2 I_d)$).

ELE SQA

Proposal Construction Summary

- Get a chain subsample $z = \{z_i\}_{i=1}^n$
- **2** Construct an RKHS sample $f \sim \mathcal{N}(\phi(x_t), \nu^2 C_z)$
- Propose x^* such that $\phi(x^*)$ is close to f (with an additional exploration term $\xi \sim \mathcal{N}(0, \gamma^2 I_d)$).

Integrate out RKHS samples f, gradient step, and ξ to obtain marginal Gaussian proposal on the input space:

$$q_{\mathbf{z}}(x^*|x_t) = \mathcal{N}(x_t, \gamma^2 I_d + \nu^2 M_{\mathbf{z}, x_t} H M_{\mathbf{z}, x_t}^{\top})$$

 $M_{\mathbf{z},x_t} = 2 \left[\nabla_x k(x, z_1) |_{x=x_t}, \dots, \nabla_x k(x, z_n) |_{x=x_t} \right],$ $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}.$

▲ 注 > < 注 > 三 二 の Q ペ 12/12/2015 12 / 24

Examples of Covariance Structure for Standard Kernels



Kameleon proposals capture local covariance structure Gaussian kernel: $k(x, x') = \exp\left(-\frac{1}{2}\sigma^{-2} ||x - x'||_2^2\right)$

$$\left[\operatorname{cov}[q_{\mathbf{z}(\cdot|\mathbf{y})}]\right]_{ij} = \gamma^2 \delta_{ij} + \frac{4\nu^2}{\sigma^4} \sum_{a=1}^n \left[k(\mathbf{y}, \mathbf{z}_a)\right]^2 (\mathbf{z}_{a,i} - \mathbf{y}_i) (\mathbf{z}_{a,j} - \mathbf{y}_j) + \mathcal{O}\left(\frac{1}{n}\right).$$

Arthur Gretton (Gatsby Unit, UCL)

Kernel Adaptive Metropolis-Hastings

12/12/2015 13 / 24

Kernel Adaptive Hamiltonian Monte Carlo (KMC)



Heiko Strathmann, Dino Sejdinovic, Samuel Livingstone, Zoltan Szabo, and Arthur Gretton, NIPS 2015

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 14 / 24

Hamiltonian Monte Carlo

- HMC: distant moves, high acceptance probability.
- Potential energy $U(q) = -\log \pi(q)$, auxiliary momentum $p \sim \exp(-K(p))$, simulate for $t \in \mathbb{R}$ along Hamiltonian flow of H(p,q) = K(p) + U(q), using operator

$$\frac{\partial K}{\partial p}\frac{\partial}{\partial q} - \frac{\partial U}{\partial q}\frac{\partial}{\partial p}$$

- 0 d¹− -2 -3-4-5-5-4-3-1 θ_2
- Numerical simulation (i.e. leapfrog) depends on gradient information.

What if gradient *unavailable*, e.g. in Bayesian GP classification?

Arthur Gretton (Gatsby Unit, UCL)

Kernel Adaptive Metropolis-Hastings

■ * * ■ * ■ = 少への 12/12/2015 15 / 24

Infinite dimensional exponential families

Proposal is RKHS exponential family model [Fukumizu, 2009; Sriperumbudur et al. 2014], but accept using true Hamiltonian (to correct for both model and leapfrog)

$$\operatorname{const} \times \pi(x) \approx \exp\left(\langle f, k(x, \cdot) \rangle_{\mathcal{H}} - A(f)\right)$$

- Sufficient statistics: feature map $k(\cdot, x) \in \mathcal{H}$, satisfies $f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{H}}$ for any $f \in \mathcal{H}$.
- Natural parameters: $f \in \mathcal{H}$.

The model is

- dense in continuous densities on compact domains (TV, KL, etc.),
- relatively robust to increasing dimensions, as opposed to e.g. KDE.

How to learn f from samples without access to A(f)?

Score matching

- Estimation of unnormalised density models from samples [Sriperumbudur et al. 2014]
- Minimises Fisher divergence

$$J(f) = \frac{1}{2} \int \pi(x) \left\| \nabla f(x) - \nabla \log \pi(x) \right\|_2^2 dx$$

Possible without accessing ∇ log π(x) and accessing π(x) only through samples x := {x_i}^t_{i=1}

$$\hat{J}(f) = \widehat{\mathbb{E}}_{x} \left\{ \sum_{\ell=1}^{d} \left[\frac{\partial^{2} f(x)}{\partial x_{\ell}^{2}} + \frac{1}{2} \left(\frac{\partial f(x)}{\partial x_{\ell}} \right)^{2} \right] \right\}$$

Expensive: full solution requires solving (td + 1)-dimensional linear system.

> < = > < = > = = < < < >

Approximate solution: KMC finite

$$f(\mathbf{x}) = \theta^\top \phi_{\mathbf{x}}$$

- Random Fourier Features $\phi_x^{\top} \phi_y \approx k(x, y)$
- $heta \in \mathbb{R}^m$ can be computed from

$$\hat{\theta}_{\lambda} := (C + \lambda I)^{-1} b$$

$$\begin{split} b &:= -\frac{1}{t} \sum_{i=1}^{t} \sum_{\ell=1}^{d} \ddot{\phi}_{x_{i}}^{\ell} \quad C := \frac{1}{t} \sum_{i=1}^{t} \sum_{\ell=1}^{d} \dot{\phi}_{x_{i}}^{\ell} \left(\dot{\phi}_{x_{i}}^{\ell} \right)^{T} \\ \text{where } \dot{\phi}_{x}^{\ell} &:= \frac{\partial}{\partial x_{\ell}} \phi_{x} \text{ and } \ddot{\phi}_{x}^{\ell} := \frac{\partial^{2}}{\partial x_{\ell}^{2}} \phi_{x}. \end{split}$$
 $\bullet \quad On-line \ updates \ cost \ \mathcal{O}(dm^{2}). \end{split}$

Updates fast, uses *all* Markov chain history. Caveat: need to initialise correctly.

Gradient norm: <u>Gaussian</u>





ELE NOR

Approximate solution: KMC lite

$$f(x) = \sum_{i=1}^{n} \alpha_i k(z_i, x)$$

- $z \subseteq x$ sub-sample.
- α from linear system

$$\hat{\alpha}_{\lambda} = -\frac{\sigma}{2}(C + \lambda I)^{-1}b$$

where $C \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n$ depend on kernel matrix

• Cost $\mathcal{O}(n^3 + n^2d)$ (or cheaper with low-rank approx., conjugate gradient).

Geometrically ergodic on logconcave targets (fast convergence).





EL SQC

19 / 24

12/12/2015

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

Does kernel HMC work in high dimensions?

Challenging Gaussian target (top):

- Eigenvalues: $\lambda_i \sim \text{Exp}(1)$.
- Covariance: diag(λ₁,...,λ_d), randomly rotate.
- Use Rational Quadratic kernel to account for resulting highly 'non-singular' length-scales.
- KMC scales up to $d \approx 30$.

An easy, isotropic Gaussian target (bottom):

• More smoothness allows KMC to scale up to $d \approx 100$.





Synthetic targets: Banana

Banana: $\mathcal{B}(b, v)$: take $X \sim \mathcal{N}(0, \Sigma)$ with $\Sigma = \text{diag}(v, 1, \dots, 1)$, and set $Y_2 = X_2 + b(X_1^2 - v)$, and $Y_i = X_i$ for $i \neq 2$. (Haario et al, 1999; 2001)



Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 21 / 24

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三回日 うらう

Synthetic targets: Banana



KMC behaves like HMC as number n of oracle samples increases.

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 22 / 24

ъ

Gaussian Process Classification on UCI data

• Standard GPC model

 $p(\mathbf{f}, \mathbf{y}, \theta) = p(\theta)p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f})$

- where $p(\mathbf{f}|\theta)$ is a GP and with a sigmoidal likelihood $p(\mathbf{y}|\mathbf{f})$.
- Goal: sample from $p(\theta|\mathbf{y}) \propto p(\theta)p(\mathbf{y}|\theta).$
- Unbiased estimate of $\hat{p}(\mathbf{y}|\theta)$ via importance sampling.
- No access to likelihood or gradient.



Gaussian Process Classification on UCI data

Standard GPC model

$$p(\mathbf{f}, \mathbf{y}, \theta) = p(\theta)p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f})$$

where $p(\mathbf{f}|\theta)$ is a GP and with a sigmoidal likelihood $p(\mathbf{y}|\mathbf{f})$.

- Goal: sample from $p(\theta|\mathbf{y}) \propto p(\theta)p(\mathbf{y}|\theta).$
- Unbiased estimate of $\hat{p}(\mathbf{y}|\theta)$ via importance sampling.
- No access to likelihood or gradient.



Significant mixing improvements over state-of-the-art.

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 23 / 24

Conclusions

- Simple, versatile, gradient-free adaptive MCMC samplers:
- Kameleon:
 - Uses local covariance structure of the target distribution at the current chain state
- Kernel HMC
 - Derivative of log density fit to samples, use this as proposal in HMC.
- Outperforms existing adaptive approaches on nonlinear target distributions
- Future work: For Kameleon, does feature space covariance track high density regions in original space? For kernel HMC, how does convergence rate degrade with increasing dimension?

• Kameleon code: https://github.com/karlnapf/kameleon-mcmc

12/12/2015

24 / 24

• Kernel HMC code: https://github.com/karlnapf/kernel hmc

 GPC model: latent process f, labels y, (with covariate matrix X), and hyperparameters θ:

$$p(\mathbf{f}, \mathbf{y}, \theta) = p(\theta)p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f})$$

where $\mathbf{f}|\theta \sim \mathcal{N}(0, \mathcal{K}_{\theta})$ is a realization of a GP with covariance \mathcal{K}_{θ} (covariance between latent processes evaluated at X).

> = = ~ ~ ~

 GPC model: latent process f, labels y, (with covariate matrix X), and hyperparameters θ:

$$p(\mathbf{f}, \mathbf{y}, \theta) = p(\theta)p(\mathbf{f}|\theta)p(\mathbf{y}|\mathbf{f})$$

where $\mathbf{f}|\theta \sim \mathcal{N}(0, \mathcal{K}_{\theta})$ is a realization of a GP with covariance \mathcal{K}_{θ} (covariance between latent processes evaluated at X).

*K*_θ: exponentiated quadratic Automatic Relevance Determination (ARD) covariance:

$$(\mathcal{K}_{\theta})_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}'_j | \theta) = \exp\left(-\frac{1}{2}\sum_{s=1}^d \frac{(x_{i,s} - x'_{j,s})^2}{\exp(\theta_s)}\right)$$

• Fully Bayesian treatment: Interested in the posterior $p(\theta|y)$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三回日 うらう

- Fully Bayesian treatment: Interested in the posterior p(heta|y)
- Cannot use a Gibbs sampler on $p(\theta, \mathbf{f}|y)$, which samples from $p(\mathbf{f}|\theta, y)$ and $p(\theta|\mathbf{f}, y)$ in turns, since $p(\theta|\mathbf{f}, y)$ is extremely sharp

- Fully Bayesian treatment: Interested in the posterior $p(\theta|y)$
- Cannot use a Gibbs sampler on $p(\theta, \mathbf{f}|y)$, which samples from $p(\mathbf{f}|\theta, y)$ and $p(\theta|\mathbf{f}, y)$ in turns, since $p(\theta|\mathbf{f}, y)$ is extremely sharp
- Filippone & Girolami, 2013 use Pseudo-Marginal MCMC to sample $p(\theta|y) = p(\theta) \int p(\theta, \mathbf{f}|y) p(\mathbf{f}|\theta) d\mathbf{f}$.

▲□ ▶ ▲ □ ▶ ▲ □ ▶ ▲□ ■ ● ● ●

- Fully Bayesian treatment: Interested in the posterior p(heta|y)
- Cannot use a Gibbs sampler on $p(\theta, \mathbf{f}|y)$, which samples from $p(\mathbf{f}|\theta, y)$ and $p(\theta|\mathbf{f}, y)$ in turns, since $p(\theta|\mathbf{f}, y)$ is extremely sharp
- Filippone & Girolami, 2013 use Pseudo-Marginal MCMC to sample $p(\theta|y) = p(\theta) \int p(\theta, \mathbf{f}|y) p(\mathbf{f}|\theta) d\mathbf{f}.$
- Unbiased estimate of $\hat{p}(\mathbf{y}|\theta)$ via importance sampling:

$$\hat{p}(\theta|\mathbf{y}) \propto p(\theta)\hat{p}(\mathbf{y}|\theta) \approx p(\theta) \frac{1}{n_{\mathrm{imp}}} \sum_{i=1}^{n_{\mathrm{imp}}} p(\mathbf{y}|\mathbf{f}^{(i)}) \frac{p(\mathbf{f}^{(i)}|\theta)}{Q(\mathbf{f}^{(i)})}$$

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

(日本)

- Fully Bayesian treatment: Interested in the posterior p(heta|y)
- Cannot use a Gibbs sampler on $p(\theta, \mathbf{f}|y)$, which samples from $p(\mathbf{f}|\theta, y)$ and $p(\theta|\mathbf{f}, y)$ in turns, since $p(\theta|\mathbf{f}, y)$ is extremely sharp
- Filippone & Girolami, 2013 use Pseudo-Marginal MCMC to sample $p(\theta|y) = p(\theta) \int p(\theta, \mathbf{f}|y) p(\mathbf{f}|\theta) d\mathbf{f}.$
- Unbiased estimate of $\hat{p}(\mathbf{y}|\theta)$ via importance sampling:

$$\hat{p}(\theta|\mathbf{y}) \propto p(\theta)\hat{p}(\mathbf{y}|\theta) \approx p(\theta) \frac{1}{n_{\mathrm{imp}}} \sum_{i=1}^{n_{\mathrm{imp}}} p(\mathbf{y}|\mathbf{f}^{(i)}) \frac{p(\mathbf{f}^{(i)}|\theta)}{Q(\mathbf{f}^{(i)})}$$

• No access to likelihood, gradient, or Hessian of the target.

(日本)

RKHS and Kernel Embedding

 For any positive semidefinite function k, there is a unique RKHS H_k. Can consider x → k(·, x) as a feature map.

RKHS and Kernel Embedding

 For any positive semidefinite function k, there is a unique RKHS H_k. Can consider x → k(·, x) as a feature map.

Definition (Kernel embedding)

Let k be a kernel on \mathcal{X} , and P a probability measure on \mathcal{X} . The kernel embedding of P into the RKHS \mathcal{H}_k is $\mu_k(P) \in \mathcal{H}_k$ such that $\mathbb{E}_P f(X) = \langle f, \mu_k(P) \rangle_{\mathcal{H}_k}$ for all $f \in \mathcal{H}_k$.

- Alternatively, can be defined by the Bochner integral $\mu_k(P) = \int k(\cdot, x) dP(x)$ (expected canonical feature)
- For many kernels k, including the Gaussian, Laplacian and inverse multi-quadratics, the kernel embedding P → μP is injective: characteristic (Sriperumbudur et al, 2010),
- captures all moments (similarly to the characteristic function).

EL SOCO

Covariance operator

Definition

The covariance operator of P is $C_P : \mathcal{H}_k \to \mathcal{H}_k$ such that $\forall f, g \in \mathcal{H}_k$, $\langle f, C_P g \rangle_{\mathcal{H}_k} = \operatorname{Cov}_P [f(X)g(X)].$

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

(日本)

Covariance operator

Definition

The covariance operator of P is $C_P : \mathcal{H}_k \to \mathcal{H}_k$ such that $\forall f, g \in \mathcal{H}_k$, $\langle f, C_P g \rangle_{\mathcal{H}_k} = \operatorname{Cov}_P [f(X)g(X)].$

- Covariance operator: C_P : H_k → H_k is given by C_P = ∫ k(·,x) ⊗ k(·,x) dP(x) − μ_P ⊗ μ_P (covariance of canonical features)
- Empirical versions of embedding and the covariance operator:

$$\mu_{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^{n} k(\cdot, z_i) \qquad \qquad C_{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^{n} k(\cdot, z_i) \otimes k(\cdot, z_i) - \mu_{\mathbf{z}} \otimes \mu_{\mathbf{z}}$$

28 / 24

12/12/2015

The empirical covariance captures **non-linear** features of the underlying distribution, e.g. Kernel PCA

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

Kernel distance gradient

$$g(x) = k(x, x) - 2k(x, y) - 2\sum_{i=1}^{n} \beta_i [k(x, z_i) - \mu_z(x)]$$
$$\nabla_x g(x)|_{x=y} = \underbrace{\nabla_x k(x, x)|_{x=y} - 2\nabla_x k(x, y)|_{x=y}}_{=0} - M_{z,y} H\beta$$

where $M_{\mathbf{z},y} = 2\left[\nabla_x k(x,z_1)|_{x=y}, \ldots, \nabla_x k(x,z_n)|_{x=y}\right]$ and $H = I_n - \frac{1}{n} \mathbf{1}_{n \times n}$

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

Cost function g



g varies most along the high density regions of the target

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 30 / 24

ELE SQA

Synthetic targets: Flower

Flower: $\mathcal{F}(r_0, A, \omega, \sigma)$, a *d*-dimensional target with:

$$\mathcal{F}(x; r_0, A, \omega, \sigma) \propto \ \exp\left(-rac{\sqrt{x_1^2 + x_2^2} - r_0 - A\cos\left(\omega atan2\left(x_2, x_1
ight)
ight)}{2\sigma^2}
ight) \times \prod_{j=3}^d \mathcal{N}(x_j; 0, 1).$$

Concentrates on r_0 -circle with a periodic perturbation (with amplitude A and frequency ω) in the first two dimensions.



EL SOCO

Synthetic targets: convergence statistics



8-dimensional $\mathcal{F}(10, 6, 6, 1)$ target; iterations: 120000, burn-in: 60000

Arthur Gretton (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

12/12/2015 32 / 24

315