Kernel Adaptive Metropolis-Hastings

Dino Sejdinovic*, Heiko Strathmann*, Maria Lomeli Garcia*, Christophe Andrieu[‡], and Arthur Gretton*

> *Gatsby Unit, CSML, University College London, [‡]School of Mathematics, University of Bristol

> > 22 June 2014



international conference on machine learning, 2014

Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

22/06/2014 1 / 15

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

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

• What proposal $q(\cdot|x_t)$?

A ∃ ► A ∃ ► ∃ ⊨ 900

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 $q(\cdot|x_t)$?
 - $\bullet\,$ Too narrow: small increments $\rightarrow\,$ slow convergence
 - Too broad: many rejections \rightarrow slow convergence

A 回 > A 回 > A 回 > 回 回 の Q (A)

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,



ELE NOR

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,



22/06/2014 3 / 15

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

EL SQA

3 / 15

22/06/2014

Sejdinovic et al (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).

EL SOCO

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

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).

4 回 > 4 回 > 4 回 > 回 回 の Q の

When is target not computable?

• Posterior inference, latent process f

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

ELE NOR

When is target not computable?

• Posterior inference, latent process f

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

• Cannot integrate out f: e.g. Gaussian process classification, θ lengthscales of covariance. MH ratio:

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

ELE SQA

4 1 1 4 1 1 1

When is target not computable?

• Posterior inference, latent process f

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

• Cannot integrate out f: e.g. Gaussian process classification, θ lengthscales of covariance. MH ratio:

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

• Replace $p(\mathbf{y}|\theta)$ with Monte Carlo estimate $\hat{p}(\mathbf{y}|\theta)$

When is target not computable?

• Posterior inference, latent process f

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

• Cannot integrate out f: e.g. Gaussian process classification, θ lengthscales of covariance. MH ratio:

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

- Replace $p(\mathbf{y}|\theta)$ with Monte Carlo estimate $\hat{p}(\mathbf{y}|\theta)$
- Replacing marginal likelihood with *unbiased estimate* still results in correct invariant distribution (Beaumont, 2003; Andrieu & Roberts, 2009)

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

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?

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

Input space \mathcal{X}

 \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 $\mathcal H$



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



Proposal Construction Summary

- Get a chain subsample $\mathbf{z} = \{z_i\}_{i=1}^n$
- **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 $\xi \sim \mathcal{N}(0, \gamma^2 I_d)$).

> = = ~ ~ ~

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 $\xi \sim \mathcal{N}(0, \gamma^2 I_d)$).

This gives:

$$x^{*}|x_{t}, f, \xi = x_{t} - \eta \nabla_{x} \|\phi(x) - f\|_{\mathcal{H}}^{2}|_{x=x_{t}} + \xi$$

JI SAR

Proposal Construction Summary

• Get a chain subsample
$$\mathbf{z} = \{z_i\}_{i=1}^n$$

- **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 $\xi \sim \mathcal{N}(0, \gamma^2 I_d)$).

This gives:

$$x^{*}|x_{t}, f, \xi = x_{t} - \eta \nabla_{x} \|\phi(x) - f\|_{\mathcal{H}}^{2} \|_{x=x_{t}} + \xi$$

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

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

$$M_{\mathbf{z},\mathbf{x}_t} = 2 \left[\nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{z}_1) |_{\mathbf{x}=\mathbf{x}_t}, \dots, \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{z}_n) |_{\mathbf{x}=\mathbf{x}_t} \right],$$

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}.$$

Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

22/06/2014 8 / 15

MCMC Kameleon

Input: unnormalized target π ; subsample size n; scaling parameters ν, γ , kernel k; update schedule $\{p_t\}_{t\geq 1}$ with $p_t \rightarrow 0$, $\sum_{t=1}^{\infty} p_t = \infty$



At iteration t + 1,

- With probability p_t, update a random subsample z = {z_i}ⁿ_{i=1} of the chain history {x_i}^{t-1}_{i=0},
- Sample proposed point x^* from $q_{\mathbf{z}}(\cdot|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t, \gamma^2 I_d + \nu^2 M_{\mathbf{z}, \mathbf{x}_t} H M_{\mathbf{z}, \mathbf{x}_t}^{\top}),$
- Accept/Reject with standard MH ratio:

$$x_{t+1} = \begin{cases} x^*, & \text{w.p. } \min\left\{1, \frac{\pi(x^*)q_{\mathsf{z}}(x_t|x^*)}{\pi(x_t)q_{\mathsf{z}}(x^*|x_t)}\right\},\\ x_t, & \text{otherwise.} \end{cases}$$

MCMC Kameleon

Input: unnormalized target π ; subsample size *n*; scaling parameters ν, γ , kernel k; update schedule $\{p_t\}_{t>1}$ with $p_t \rightarrow$ 0. $\sum_{t=1}^{\infty} p_t = \infty$



At iteration t + 1.

- With probability p_t , update a random subsample $\mathbf{z} = \{z_i\}_{i=1}^n$ of the chain history $\{x_i\}_{i=0}^{t-1}$,
- 2 Sample proposed point x^* from $q_{\mathbf{z}}(\cdot|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t, \gamma^2 I_d + \nu^2 M_{\mathbf{z}, \mathbf{x}_t} H M_{\mathbf{z}, \mathbf{x}_t}^{\top}),$

Accept/Reject with standard MH ratio:

$$x_{t+1} = \begin{cases} x^*, & \text{w.p. } \min\left\{1, \frac{\pi(x^*)q_z(x_t|x^*)}{\pi(x_t)q_z(x^*|x_t)}\right\},\\ x_t, & \text{otherwise.} \end{cases}$$

Convergence to target π preserved as long as $p_t \to 0$ (Roberts & Rosenthal, 2007). > < = > < = > = = = < < < < Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

Locally aligned covariance



Kameleon proposals capture local covariance structure

Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

22/06/2014 10 / 15

ELE NOR

Locally aligned covariance



Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

22/06/2014 10 / 15

EL SQA

Examples of Covariance Structure for Standard Kernels

• Linear kernel:
$$k(x, x') = x^{\top}x'$$

$$q_{\mathsf{z}}(\cdot|y) = \mathcal{N}(y, \gamma^2 I + 4\nu^2 \mathsf{Z}^\top \mathsf{H} \mathsf{Z})$$

classical Adaptive Metropolis Haario et al 1999;2001.

ELE SQC

Examples of Covariance Structure for Standard Kernels

• Linear kernel:
$$k(x, x') = x^{\top}x'$$

$$q_{\mathsf{z}}(\cdot|y) = \mathcal{N}(y, \gamma^2 I + 4\nu^2 \mathsf{Z}^\top \mathsf{H} \mathsf{Z})$$

classical Adaptive Metropolis Haario et al 1999;2001.

• Gaussian kernel: $k(x, x') = \exp\left(-\frac{1}{2}\sigma^{-2} ||x - x'||_{2}^{2}\right)$

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

Influence of previous points z_a on covariance is weighted by similarity $k(y, z_a)$ to current location y.

22/06/2014

11 / 15

UCI Glass dataset



comparison in terms of all mixed moments up to order 3

8-dimensional non-linear posterior $p(\theta|\mathbf{y})$: no ground truth, performance with respect to a long-run, heavily thinned benchmark sample.

ELE DOO

12 / 15

UCI Glass dataset



comparison in terms of all mixed moments up to order 3

8-dimensional non-linear posterior $p(\theta|\mathbf{y})$: no ground truth, performance with respect to a long-run, heavily thinned benchmark sample.

ELE DOO

12 / 15

UCI Glass dataset



comparison in terms of all mixed moments up to order 3

8-dimensional non-linear posterior $p(\theta|\mathbf{y})$: no ground truth, performance with respect to a long-run, heavily thinned benchmark sample.

ELE NOR

12 / 15

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)



Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

22/06/2014 13 / 15

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三回日 うらう

Synthetic targets: convergence statistics



Strongly twisted 8-dimensional $\mathcal{B}(0.1, 100)$ target; iterations: 80000, burn-in: 40000

Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

Conclusions

- A simple, versatile, gradient-free adaptive MCMC sampler
- Proposals automatically conform to the local covariance structure of the target distribution at the current chain state
- Outperforms existing approaches on nonlinear target distributions
- Future directions: tradeoff between the sub-sampling and convergence; samplers on non-Euclidean domains

ocode: https://github.com/karlnapf/kameleon-mcmc

EL SOCO

 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).

ELE DOG

 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:

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

 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:

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

• $p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^{n} p(y_i|f_i)$ is a product of sigmoidal functions:

$$p(y_i|f_i) = \frac{1}{1 - \exp(-y_i f_i)}, \qquad y_i \in \{-1, 1\}.$$

• 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 $p(\theta|y)$
- 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

▲□ ▶ ▲ □ ▶ ▲ □ ▶ □ □ ● ○ ○ ○

- 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 $p(\theta|y)$
- 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 $p(\theta|y) = p(\theta) \int p(\theta, \mathbf{f}|y) p(\mathbf{f}|\theta) d\mathbf{f}$.
- Unbiased estimate of $\hat{p}(\mathbf{y}|\theta)$ via importance sampling:

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

▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ― 臣 ⊨ → のへ(?)

- 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 $p(\theta|y) = p(\theta) \int p(\theta, \mathbf{f}|y) p(\mathbf{f}|\theta) d\mathbf{f}$.
- Unbiased estimate of $\hat{p}(\mathbf{y}|\theta)$ via importance sampling:

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

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

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.

∃ >

JI SAR

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 $P \mapsto \mu_P$ 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 $\forall f, g \in \mathcal{H}_k$, $\langle f, C_P g \rangle_{\mathcal{H}_k} = \operatorname{Cov}_P [f(X)g(X)].$

Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

▲母 ▲ ヨ ▲ ヨ ★ ヨ ヨ ろ Q ()

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:

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

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

19 / 15

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

22/06/2014 20 / 15

Cost function g



\boldsymbol{g} varies most along the high density regions of the target

Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

22/06/2014 21 / 15

<ロト < 聞 > < 臣 > < 臣 > 三日 の Q @

Synthetic targets: Flower

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

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

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



JIN NOR

Synthetic targets: convergence statistics



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

Sejdinovic et al (Gatsby Unit, UCL) Kernel Adaptive Metropolis-Hastings

22/06/2014 23 / 15

리님