

Marginal MCMC with Noisy Local Approximations

¹**Andrew D. Davis**

²Youssef Marzouk, ³Aaron Smith, ⁴Natesh Pillai

¹Cold Region Research and Engineering Labs (CRREL)

²Massachusetts Institute of Technology

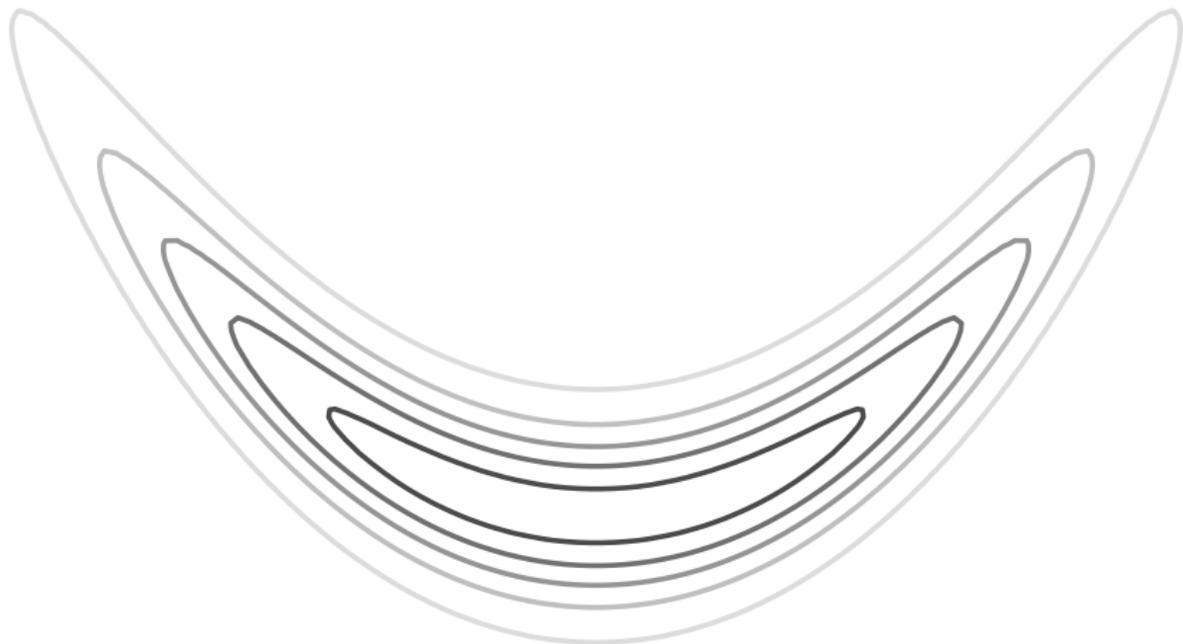
³University of Ottawa

⁴Harvard University

February 27, 2019

Problem statement

We want to **characterize the distribution π** using MCMC



We want to **characterize the distribution π** using MCMC

Two problems:

- 1 The probability density function $\pi(x)$ is **computationally expensive**



Problem statement

We want to **characterize the distribution π** using MCMC

Two problems:

- 1 The probability density function $\pi(x)$ is **computationally expensive**
- 2 Only **noisy** density evaluations are available $\hat{\pi}(x) \approx \pi(x)$



Bayesian inference

- ▶ Quantity of interest X (random variable)
- ▶ Data Y (random variable)

Bayes' theorem updates our degree of certainty given:

- ① Data/observations
- ② Physical models

Bayesian inference

- ▶ Quantity of interest X (random variable)
- ▶ Data Y (random variable)

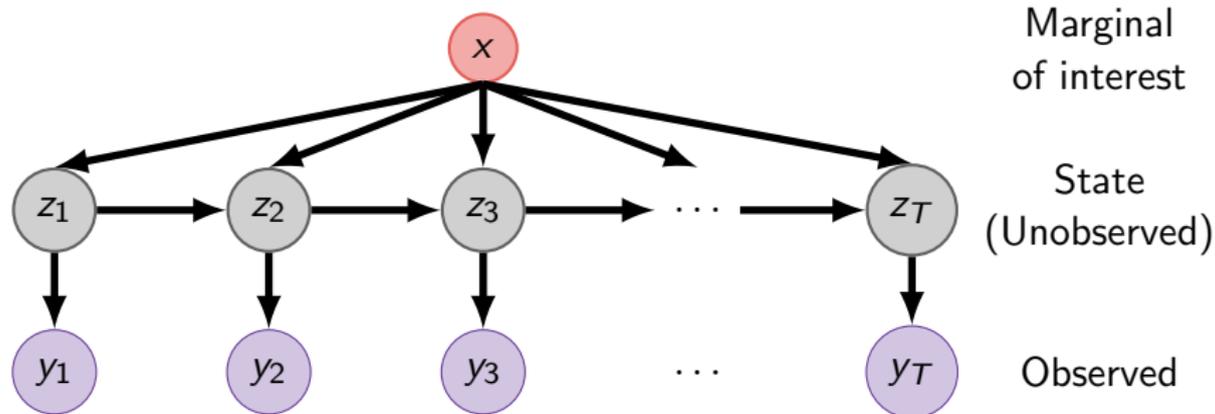
Bayes' theorem updates our degree of certainty given:

- ① Data/observations
 - ② Physical models
- ▶ **Prior** distribution: π_X models our degree of certainty about X
 - ▶ **Likelihood** distribution: $\pi_{Y|X}$ quantifies model-data mismatch
 - ▶ **Posterior** distribution: $\pi_{X|Y}$ updates the prior given Y
 - ▶ **Bayes' rule** relates the densities using the law of total probability

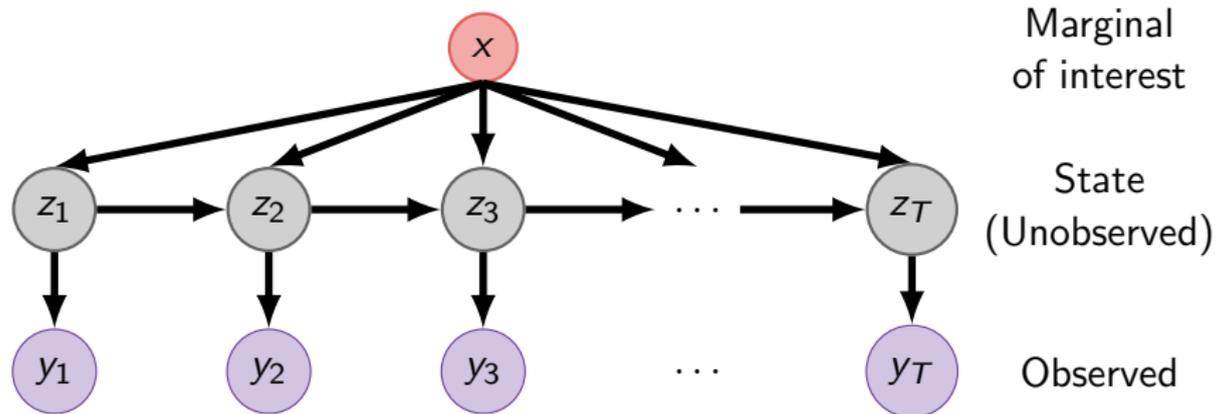
$$\pi(x|y) \propto \pi(x)\pi(y|x)$$

Static parameter estimation—

- ▶ Given: data $y_{1:T}$ with an unobservable state $z_{1:T}$
- ▶ We want to characterize the distribution over static parameters $x|y_{1:T}$
 - ▶ x is low-dimensional



Marginalize state variables



Marginalizing avoids characterizing the joint density over $[x, z_{1:T}]$

$$\underbrace{\pi(x|y_{1:T})}_{\text{Posterior}} \propto \underbrace{\pi(x)}_{\text{Prior}} \underbrace{\int \pi(z_{1:T}, y_{1:T}|x) dz}_{\text{Likelihood}}$$

Partition parameter space:

- ① Coordinates characterized by MCMC (\mathbf{x})
 - ② Coordinates to be “marginalized away” ($z_{1:T}$)
- ▶ Model evaluations are (even more) computationally expensive
 - ▶ We only have noisy estimates of the likelihood

$$\hat{\pi}(y_{1:T}|\mathbf{x}) \approx \int \pi(z_{1:T}, y_{1:T}|\mathbf{x}) dz$$

- ▶ Pseudo-marginal MCMC characterizes marginal distributions
Beaumont, 2010, Andrieu and Roberts, 2009

- ▶ Exploit the target density's regularity to build a surrogate model
 - ▶ Polynomial approximations Marzouk et al., 2007, Marzouk and Xiu, 2009
 - ▶ Gaussian processes Rasmussen, 2006, Bernardo et al., 2008, Sacks et al., 1989, Santner et al., 2003

- ▶ Continual surrogate refinement asymptotically guarantees we characterize the true target distribution Conrad et al., 2016

- ▶ Trigger surrogate refinement using a specialized error indicator for local polynomial approximations
- ▶ Derive an ideal refinement rate using a surrogate-bias and MCMC-variance trade-off
- ▶ Build the surrogate model from noisy target density evaluations

- ▶ Given: n (potentially noisy) density evaluations at points $\{\mathbf{x}^{(i)}\}_{i=1}^n$
- ▶ **Minimize the least squares error** between $\pi(\mathbf{x})$ and a degree p polynomial $\check{\pi}(\mathbf{x})$

Conn, 2009, Stone, 1977, Kohler, 2002

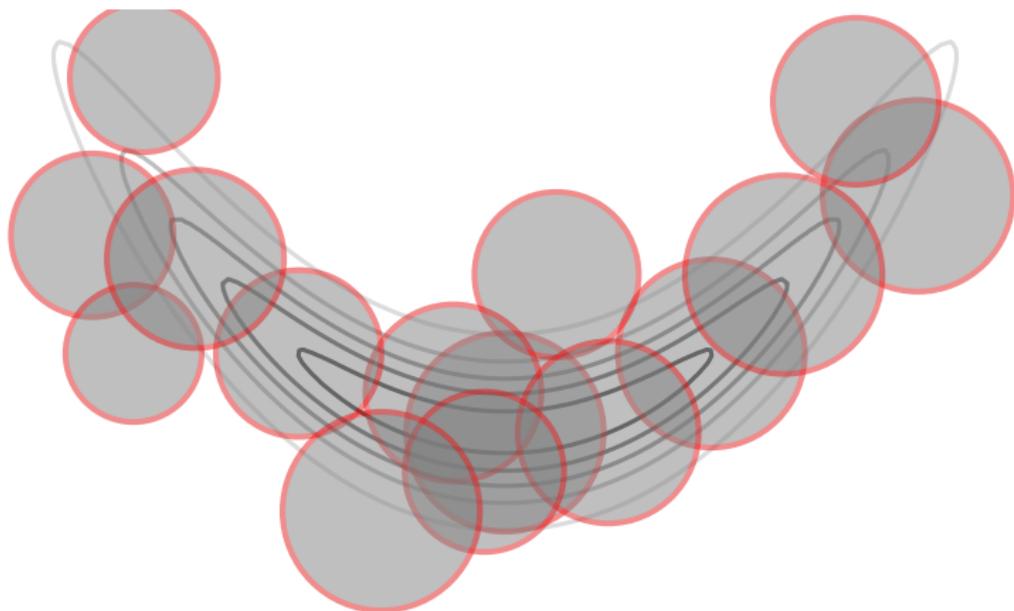
$$\check{\pi}(\mathbf{x}) = \arg \min_{\rho(\mathbf{x})} \sum_{i=1}^n \left(\rho(\mathbf{x}^{(i)}) - \pi(\mathbf{x}^{(i)}) \right)^2 K(\mathbf{x}^{(i)}, \mathbf{x})$$

- ▶ *Locally* supported kernel $K(\mathbf{x}', \mathbf{x})$ finds k nearest neighbors
 - ▶ Note: this is NOT a Markov transition kernel
- ▶ $\mathcal{B}_k(\mathbf{x})$: smallest ball centered at \mathbf{x} with k density evaluations

$$K(\mathbf{x}', \mathbf{x}) = \begin{cases} 1 & \mathbf{x}' \in \mathcal{B}_k(\mathbf{x}) \\ 0 & \text{otherwise} \end{cases}$$

Local polynomial surrogates

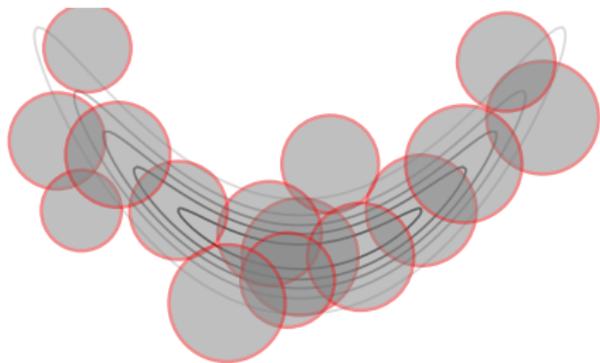
Build a local approximation in a ball around each point ...



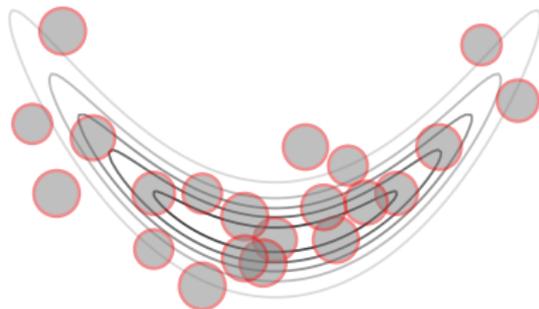
Surrogate convergence

In the **exact evaluation** case, the surrogate $\tilde{\pi}(x)$ approaches $\pi(x)$ as:

- 1 The ball size $\Delta \rightarrow 0$ (number of evaluations $n \rightarrow \infty$)



Large balls

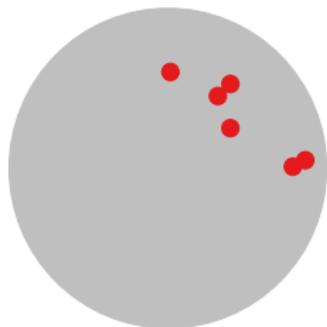


Small balls

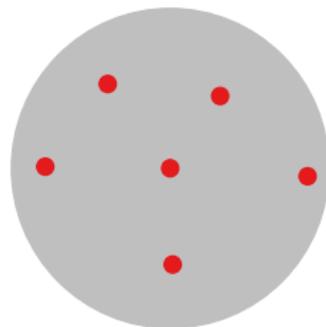
Surrogate convergence

In the **exact evaluation** case, the surrogate $\tilde{\pi}(x)$ approaches $\pi(x)$ as:

- 1 The ball size $\Delta \rightarrow 0$ (number of evaluations $n \rightarrow \infty$)
- 2 Λ -poisedness is maintained inside each ball



Poorly poised



Well poised

Recall: Markov chain Monte Carlo (MCMC)

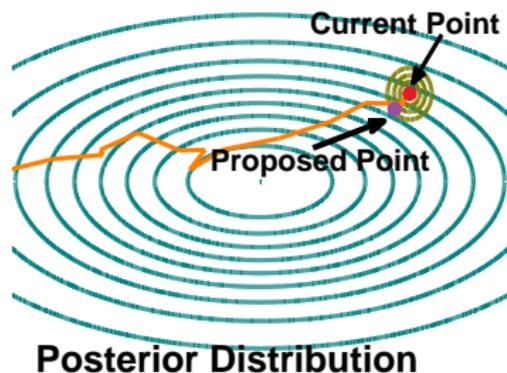
Three step algorithm:

- 1 Propose $x' \sim q_X(\cdot|x^{(t)})$
- 2 Acceptance probability

$$\alpha = \min \left(1, \frac{\pi(x')q(x^{(t)}|x')}{\pi(x^{(t)})q(x'|x^{(t)})} \right)$$

- 3 Accept/reject

$$x^{(t+1)} = \begin{cases} x' & \text{with probability } \alpha \\ x^{(t)} & \text{else} \end{cases}$$



Metropolis et al., 1953

Hastings, 1970

and variations . . .

Haario et al., 2006

Parno and Marzouk, 2014

Brooks et al., 2011

Recall: Markov chain Monte Carlo (MCMC)

Three step algorithm:

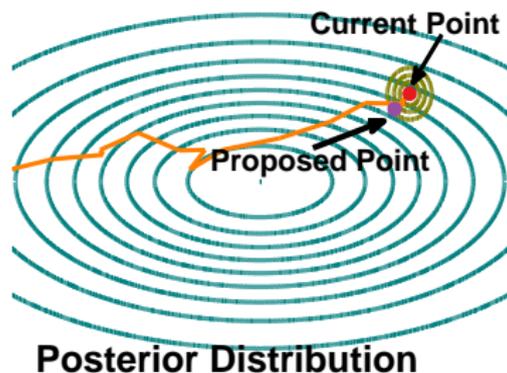
- 1 Propose $x' \sim q_X(\cdot|x^{(t)})$
- 2 Acceptance probability

$$\alpha = \min \left(1, \frac{\pi(x')q(x^{(t)}|x')}{\pi(x^{(t)})q(x'|x^{(t)})} \right)$$

- 3 Accept/reject

$$x^{(t+1)} = \begin{cases} x' & \text{with probability } \alpha \\ x^{(t)} & \text{else} \end{cases}$$

- ▶ MCMC requires an expensive density evaluation every step
- ▶ Computationally prohibitive



Metropolis et al., 1953
Hastings, 1970
and variations . . .
Haario et al., 2006
Parno and Marzouk, 2014
Brooks et al., 2011

$\check{\pi}(x) \approx \pi(x)$ is a local polynomial approximation of the target density

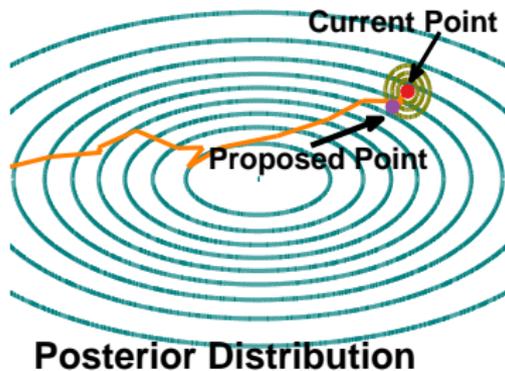
Four step algorithm:

- 1 Possibly refine $\check{\pi}(\cdot)$ near $x^{(t)}$
- 2 Propose $x' \sim q_X(\cdot|x^{(t)})$
- 3 Acceptance probability

$$\alpha = \min \left(1, \frac{\check{\pi}(x')q(x^{(t)}|x')}{\check{\pi}(x^{(t)})q(x'|x^{(t)})} \right)$$

- 4 Accept/reject

$$x^{(t+1)} = \begin{cases} x' & \text{with probability } \alpha \\ x^{(t)} & \text{else} \end{cases}$$

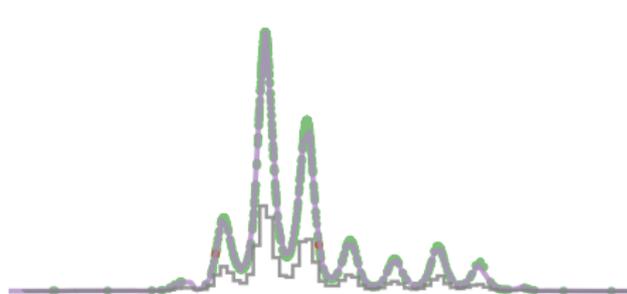


We devise a refinement strategy such that:

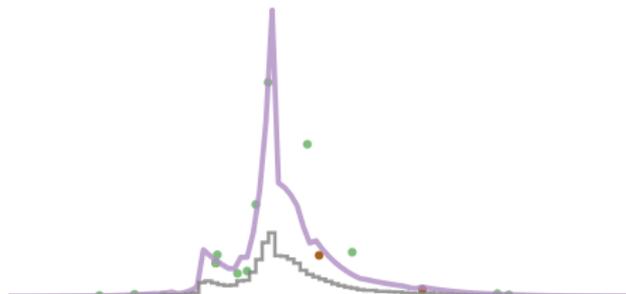
- ▶ Refinements cannot happen every MCMC step
- ▶ The expected number of refinements is infinite as $t \rightarrow \infty$
- ▶ The refinement frequency decays as $t \rightarrow \infty$
- ▶ We choose new points that maintain Λ -poisedness

Over- versus under-refinement

- ▶ Two error sources after a *finite* number of MCMC steps:
 - ① MCMC variance (CLT: decays like $1/t$)
 - ② Surrogate bias
- ▶ Frequent refinement \Rightarrow MCMC variance dominates error
- ▶ Infrequent refinement \Rightarrow surrogate bias dominates error



Over-refined



Under-refined

- ▶ Assume the MCMC variance decays with the number of steps

$$[\text{MCMC variance}] \leq C_{\text{MCMC}} t^{-1}$$

- ▶ The surrogate bias is bounded

$$[\text{Surrogate bias}] = |\hat{\pi}(\mathbf{x}) - \pi(\mathbf{x})| \leq C_{\text{Surrogate}} \bar{\Lambda} \Delta^{p+1}$$

- ▶ **Balance MCMC variance with surrogate bias squared**

$$C_{\text{MCMC}} t^{-1} \sim C_{\text{Surrogate}}^2 \bar{\Lambda}^2 \Delta^{2(p+1)}$$

Balance MCMC variance with surrogate bias squared

- ▶ Divide the chain into m levels
- ▶ Prescribe an error threshold $\gamma(m) = \gamma_0 m^{-\gamma_1}$ on each level
 - ▶ γ_0 : initial error threshold
 - ▶ γ_1 : error threshold decay rate

- ▶ Switch to level $m + 1$ when

$$C_{\text{MCMC}} t_m^{-1} = \gamma_0^2 m^{-2\gamma_1}$$

- ▶ Since C_{MCMC} is unknown

$$t_m = \varphi_1 M^{2\gamma_1}$$

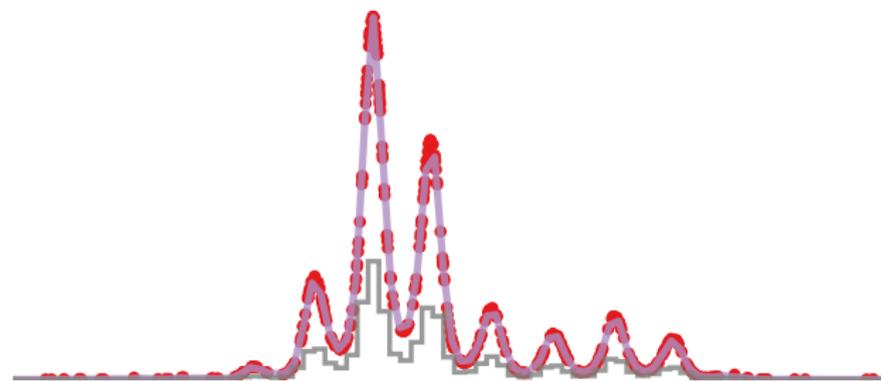
- ▶ Require $\gamma_1 > 0.5$ so the length of each level $t_m - t_{m-1}$ grows

Structural refinement

At $x^{(t)}$, refine the surrogate if:

- 1 The local error indicator is greater than the level's threshold

$$\Delta^{p+1}(x^{(t)}) > \gamma_0 M^{-\gamma_1}$$



Red dots:

Exact density
evaluations

Purple line:

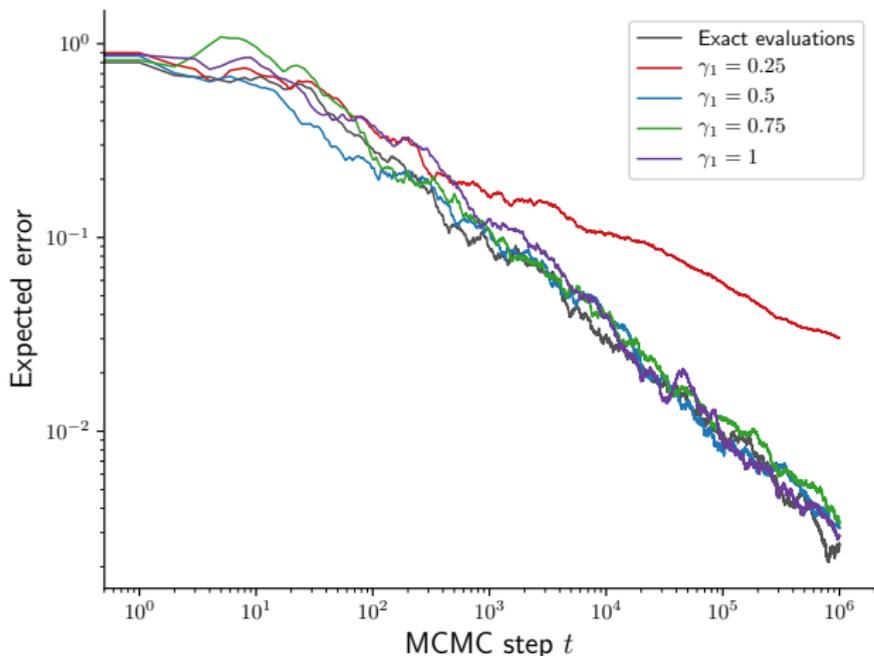
Local polynomial
approximation

Grey line:

Binned MCMC
samples

Expected error (structural refinement)

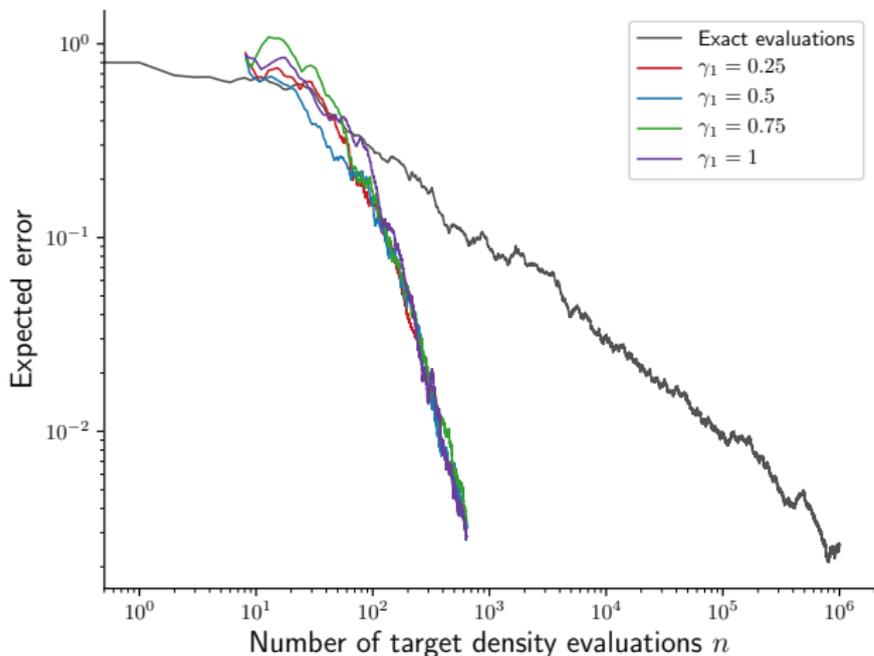
- ▶ MCMC with exact evaluations is a lower bound for the expected error
- ▶ The expected error decays at the same rate as MCMC with exact evaluations ($\gamma_1 > 0.5$)



▶ Refinement locations are chosen by the local error indicator

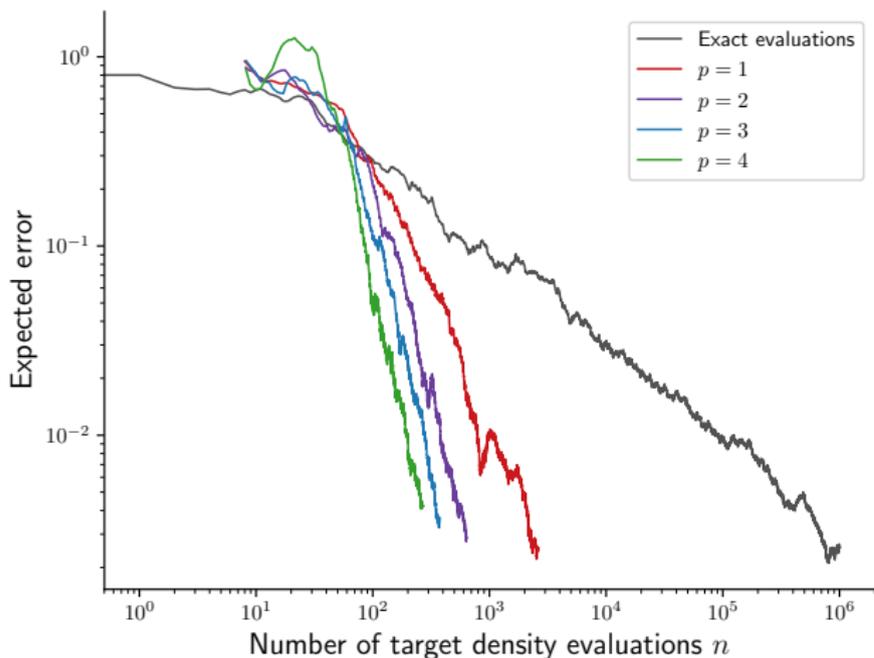
Expected error (structural refinement)

- ▶ The expected number of refinements is the same (when $\gamma_1 > 0.5$)
- ▶ When $\gamma_1 \leq 0.5$, the surrogate is underrefined
 - ▶ The error is dominated by the structural bias



Expected error (structural refinement)

- Efficiency improves with polynomial order



Local polynomial surrogates

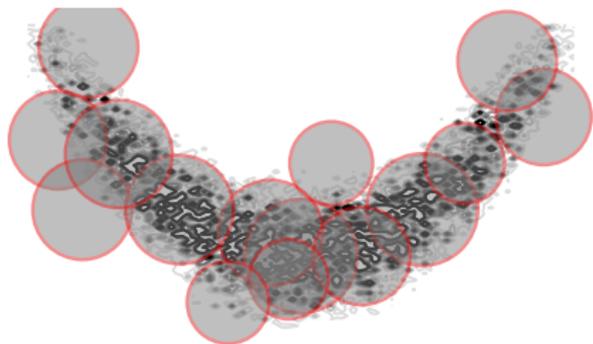
Consider situations with **noisy evaluations** of the target density



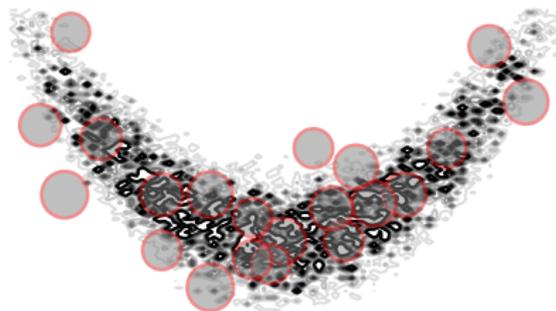
Surrogate convergence (noisy density evaluations)

The expected surrogate $\mathbb{E}[\check{\pi}(x)]$ approaches $\pi(x)$ as:

- 1 The ball size $\Delta \rightarrow 0$ (number of evaluations $n \rightarrow \infty$)
- 2 Λ -poisedness is maintained inside each ball



Large balls

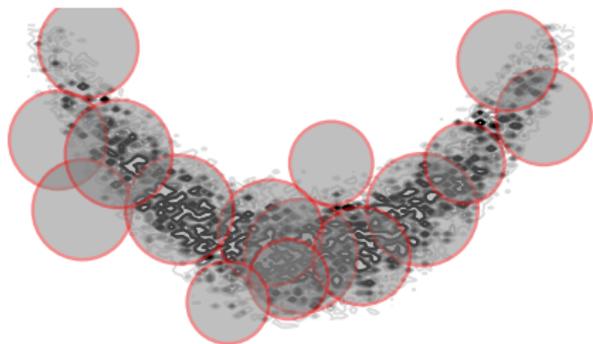


Small balls

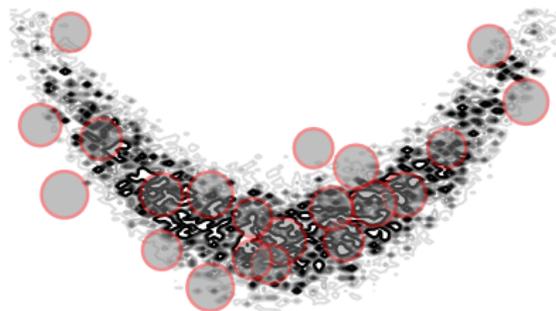
Surrogate convergence (noisy density evaluations)

The expected surrogate $\mathbb{E}[\check{\pi}(x)]$ approaches $\pi(x)$ as:

- 1 The ball size $\Delta \rightarrow 0$ (number of evaluations $n \rightarrow \infty$)
- 2 Λ -poisedness is maintained inside each ball



Large balls



Small balls

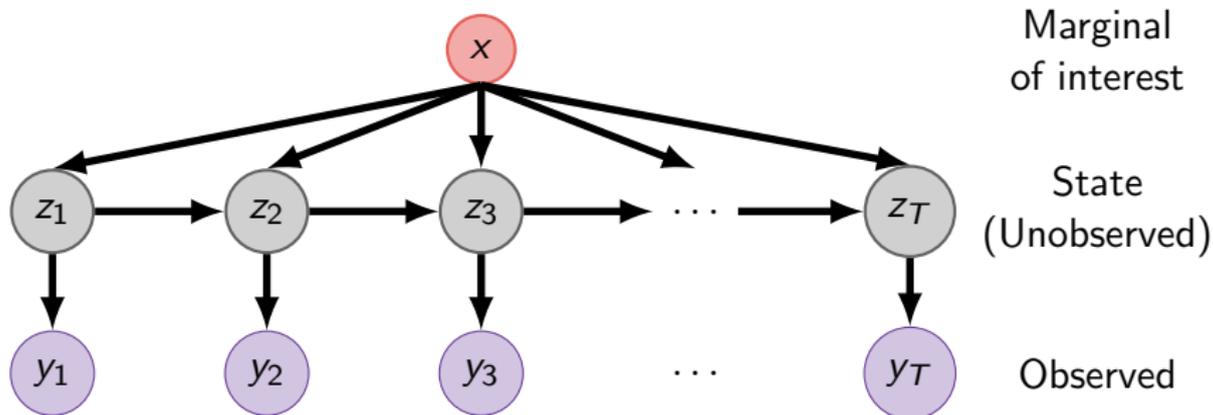
Replacing $\pi(x)$ with $\check{\pi}(x)$ within MCMC asymptotically characterizes the target distribution Davis et al., in progress

State space example: stochastic volatility

- ▶ Hyperparameters: $[\nu, \phi, \sigma^2] = g(x_1, x_2, x_3)$
- ▶ g is nonlinear transformation such that $x \sim N(0, I)$ implies $[\nu, \phi, \sigma^2]$ are sampled from Gamma distributions

$$Z_1 \sim N(\nu, (1 - \phi^2)^{-1}) \quad \text{and} \quad Z_t | Z_{t-1} \sim N(\nu + \phi(Z_{t-1} - \nu), \sigma^2)$$

$$Y_t \sim N(0, \exp(Z_t))$$



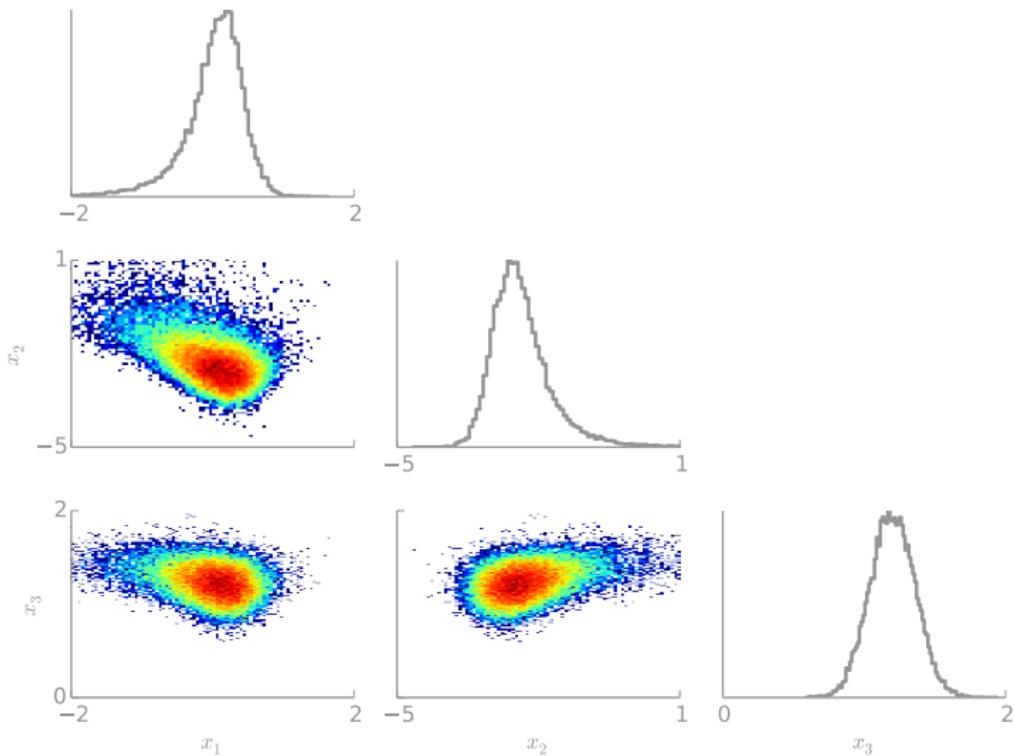
Problem statement: stochastic volatility

- ▶ Given: data $y_{1:T}$
- ▶ Define the target posterior marginal distribution

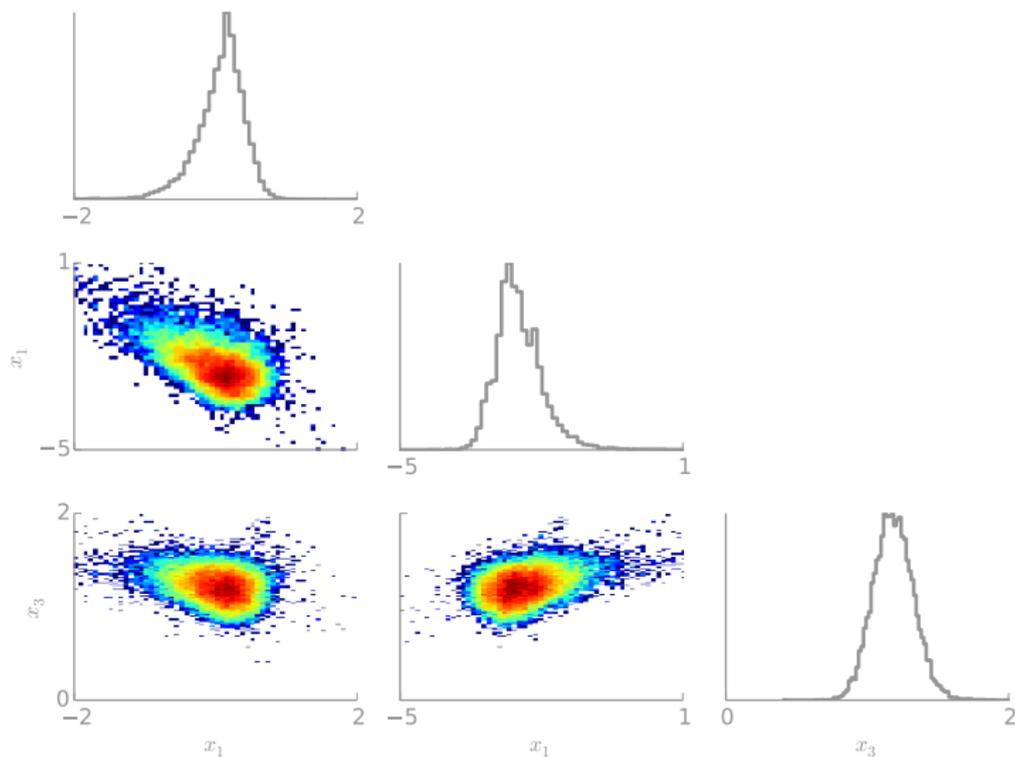
$$\begin{aligned}\pi(\mathbf{x}|y_{1:T}) &= \int \pi(\mathbf{x}, z_{1:T}|y_{1:T}) dz_{1:T} \\ &\propto \underbrace{\pi(\mathbf{x})}_{\text{Prior}} \underbrace{\int \pi(y_{1:T}|z_{1:T}, \mathbf{x})\pi(z_{1:T}|\mathbf{x}) dz_{1:T}}_{\text{Marginal likelihood}}\end{aligned}$$

- ▶ Estimate the marginal likelihood with importance sampling
 - ▶ Choosing the biasing distribution is nontrivial!

Stochastic volatility results: PM-MCMC



Stochastic volatility results: LA-MCMC



► Number of importance sampling estimates: $\approx 2.8 \times 10^3 \ll 2 \times 10^5$

- ▶ Building and refining a local polynomial approximations significantly reduces computational expense
- ▶ An ideal refinement rate comes from a bias-variance trade-off between surrogate model error and MCMC variance defines ideal
- ▶ We characterize distributions given *noisy* target density evaluations by building an asymptotically exact surrogate
 - ▶ e.g., posterior marginal distributions
- ▶ Code: MIT Uncertainty Quantification (MUQ) libraries
`muq.mit.edu`