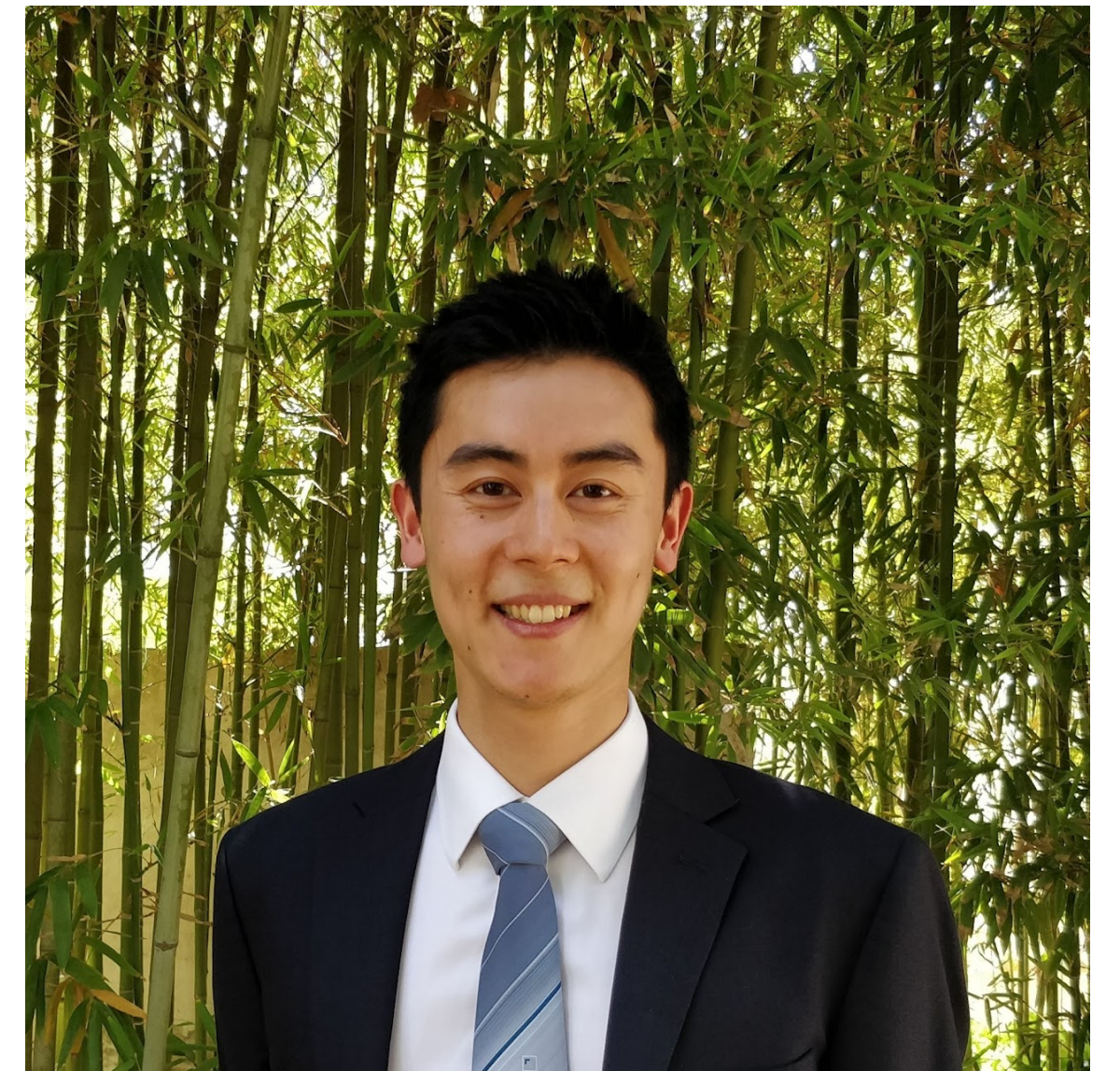


Comparison Theorems for Practical Slice Sampling

Kavli Institute for Cosmology, 5 November 2024

Sam Power, University of Bristol

joint work with: D. Rudolf (Passau), B. Sprungk (Freiberg), A.Q. Wang (Warwick)



feel free to stop me at any point

Markov chain Monte Carlo

- “target” distribution π on \mathbf{R}^d
- want samples from π to answer questions
- MCMC: use *iterative* strategy to obtain *approximate* samples

$$X_0 \rightarrow X_1 \rightarrow X_2 \rightarrow \cdots \rightarrow X_T \overset{d}{\approx} \pi$$

$$\frac{1}{T} \sum_{0 < t \leq T} f(X_t) \approx \int \pi(\mathrm{d}x) f(x) =: \pi(f)$$

Some challenges in MCMC

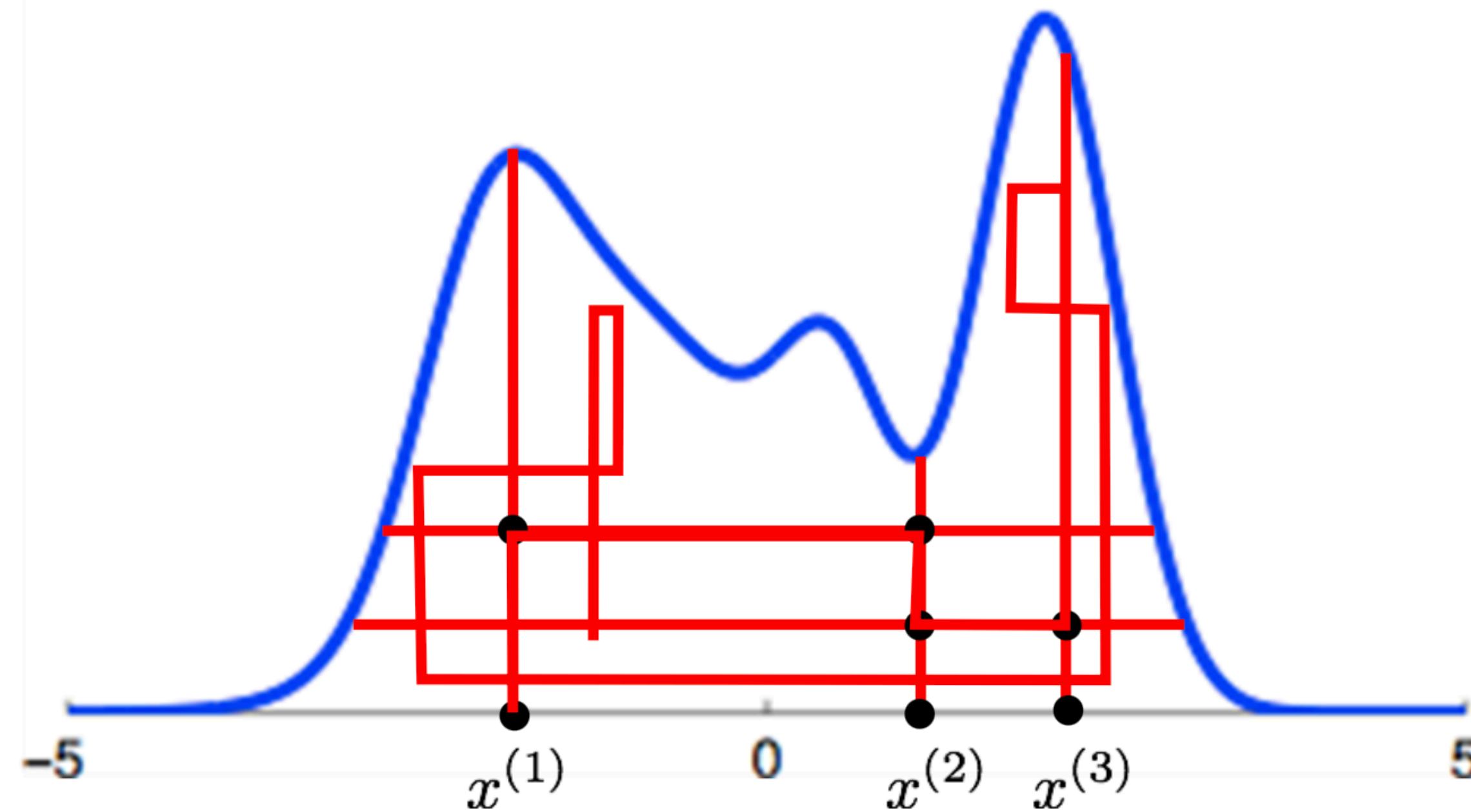
- designing effective Markov kernels
- obtaining and using useful information about π
- tuning of algorithm hyperparameters (step-size, etc.)

Slice Sampling for MCMC

- assume that we can only compute density of π (up to a constant)
- trick: sampling from π is equivalent to sampling *uniformly* under its graph
- mathematically: $\Pi(\mathrm{d}x, \mathrm{d}t) = \mathbf{1} [0 \leq t \leq \pi(x)] \mathrm{d}x \mathrm{d}t$

Slice Sampling

Define a Markov chain that samples uniformly from the area beneath the curve. This means that we need to introduce a “height” into the MCMC sampler.



(from slides of Ryan Adams)

Some Useful Definitions

- write $\bar{\pi} = \max \{ \pi(x) : x \in \mathcal{X} \}$, $T = [0, \bar{\pi}]$ ('heights')
- write $G(t) = \{x : \pi(x) \geq t\}$ for the super-level set ('slice')
- write $\nu_t = \text{Unif}(G(t))$ ('on-slice distribution')
- write $m(t) = \text{Leb}(G(t))$ ('mass function')

Slice Sampling: Algorithm

- want to generate sequence $\{ (X_n, T_n) : n \geq 1 \}$
- so,
 - given $X_{n-1} = x$, sample $T_n \sim \text{Unif} \left([0, \pi(x)] \right)$
 - (sample a height)
 - given $T_n = t$, sample $X_n \sim \nu_t$
 - (sample uniformly ‘on the slice’)

Slice Sampling: Qualitatively

- under mild conditions, gives an ergodic, π -invariant Markov chain
 - \rightsquigarrow fit for purpose in MCMC
- under still mild conditions, is even *exponentially* convergent
 - \rightsquigarrow bonus results, e.g. Markov chain CLT
- surprisingly hard to break

Slice Sampling: Dualities

- in principle, we could look at the convergence properties of
 1. $\{X_n\}$ alone (of key interest)
 2. $\{(X_n, T_n)\}$ jointly
 3. $\{T_n\}$ alone (univariate)
- in fact, all three processes will equilibrate ‘identically’

Slice Sampling: Invariances

- consider two target distributions π_1, π_2 with mass functions m_1, m_2
- if $m_1 = m_2$, then convergence profiles will be *identical*
 - follows from considering $\{T_n\}$ chain
- convergence is thus (nominally) agnostic to
 - { dimension, linear change of variables, ‘rearrangement’ of mass, ... }
- conclusion: can ‘WLOG’ very freely!
- consequence: convergence is ideally quite robust

Slice Sampling: Quantitatively

- for specific π , strong quantitative theory available
 - π spherically-symmetric, log-concave \rightsquigarrow ‘relaxation time’ $\sim \text{dim}$
 - π multivariate Student-t \rightsquigarrow ‘relaxation time’ $\sim \text{dim}^2$
 - (other explicit examples can be studied)
- noteworthy: barely slowed down by heavy tails; *rare* property

Implementing Slice Sampling

“given $T_n = t$, sample $X_n \sim \nu_t$ ”

– *Sam Power, Slide 8*

Life on the Slice

- recall $\nu_t = \text{Unif}(G(t))$, $G(t) = \{x : \pi(x) \geq t\}$
- if $G(t)$ is a { ball, box, simplex, ... }, then sampling from ν_t is fine
- if not, then we have a new problem

Hybrid Slice Sampling

- instead of
 - “given $T_n = t$, sample $X_n \sim \nu_t$ ”
- do
 - given $X_{n-1} = x, T_n = t$,
 - sample $X_n \sim \text{MCMC} (x \rightarrow x'; \text{target} = \nu_t)$
 - (call this Markov kernel H_t)

Properties of HSS

- this new algorithm ...
 - is still implementable
 - still has the right qualitative behaviour (invariance, ergodicity, ...)
 - is not as (statistically) efficient as the ‘ideal’ Slice Sampler
- how do we quantify the cost of approximation?
 - \rightsquigarrow Markov chain comparison theory

A word on MCMC theory

- it is (IMO) nice to have an idea of how long it takes for MCMC to converge
- generally, we first seek to understand complexity w.r.t. relevant features
 - { dimension, curvature, dependence, ... }
- reasonably useful to **compare** algorithms; “*which MCMC should I use?*”
- relatively rare (AFAIK) to use such results to e.g. set algorithm runtime

On Details

- technical details about the comparison framework will be suppressed
- general strategy: “ L^2 / Hilbert space / Dirichlet form” approach
- “ K converges to π at rate γ ” \sim “spectral gap” $\geq \gamma$, reversible + positive
- this notion of convergence is reasonably strong
- after digesting the definitions, the techniques are relatively **easy to use**

General Strategy

- let U = Ideal Slice Sampling, H = Hybrid Slice Sampling
- to quantify how quickly H converges to π , we will study
 - how quickly U converges to π , and
 - how well H approximates U
 - (rather, how well H_t approximates ν_t)

Realities of Comparison Theory

- to say that H gives a good Markov chain, we are arguing that
 - U gives a good Markov chain, and
 - H is a good approximation of U
- in principle, H could fail to approximate U well, but still work well
 - our analysis would fail to capture this (\exists **examples**)

Some Warm-Up Results

- in complete generality, U dominates H
 - \rightsquigarrow all else being equal, always prefer ideal chain
- for experts: any Metropolis(~~Hastings~~) kernel **can be written** as such an H
 - \rightsquigarrow all such chains are automatically dominated by (ideal) SS
- conversely: we are interested in when H is **almost as good** as U

A Generic Result

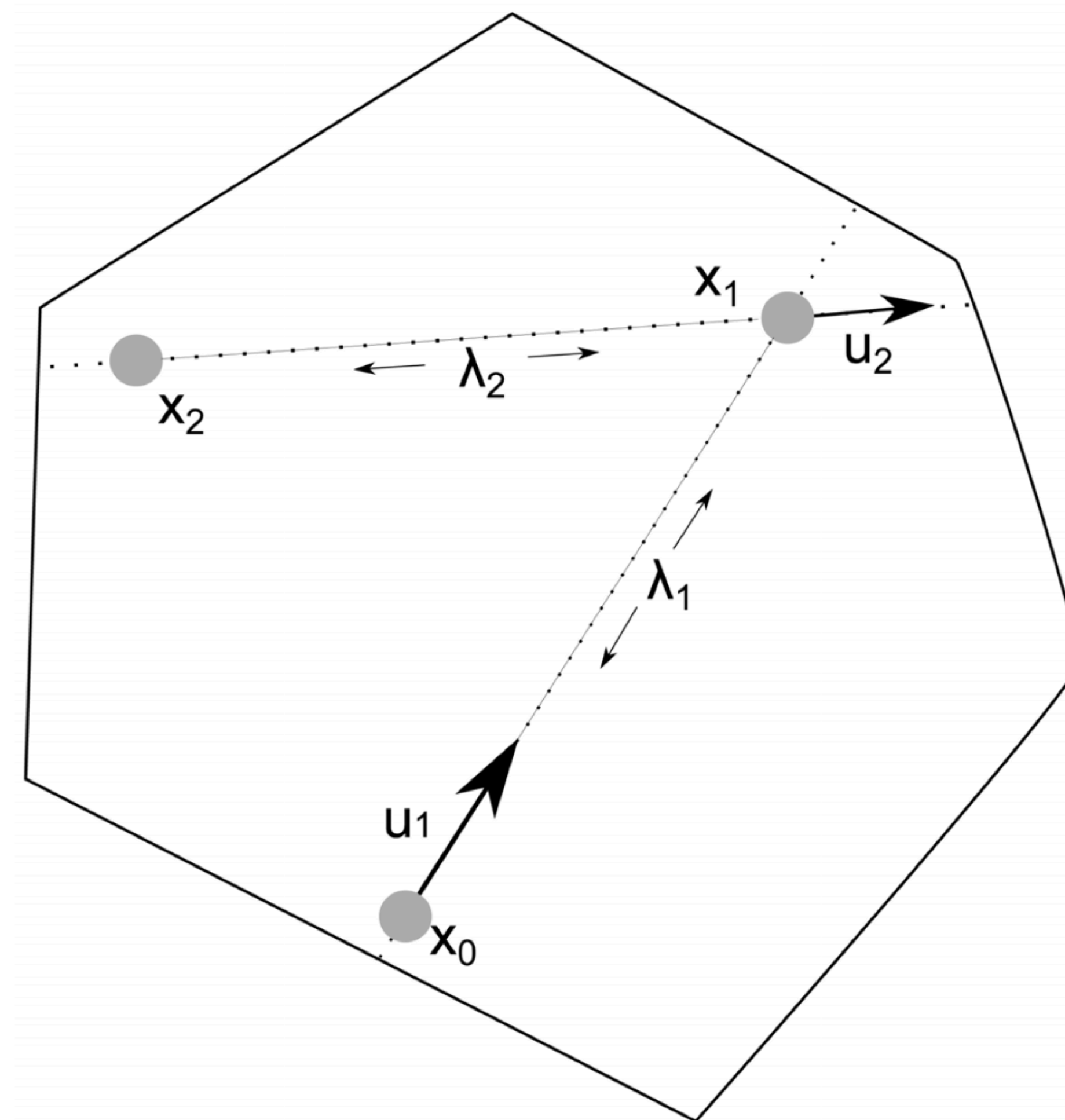
- suppose that
 - for $t \in T$, the on-slice kernel H_t converges to ν_t at rate $\sigma_t > 0$, and
 - $\sigma_H := \inf_{t \in T} \sigma_t > 0$
- the convergence rates of H and U then satisfy

$$\gamma_H \geq \sigma_H \cdot \gamma_U,$$

- interpretation: HSS is at most a factor σ_H ‘worse’ than ideal SS

Case Study: Hit-and-Run on the Slice

- simple method for sampling uniform distributions on convex body G
- at $X_{n-1} = x$,
 - sample $U_n \sim \text{Unif}(\mathbb{S}^{d-1})$
 - look at $(x + U_n \mathbf{R}) \cap G$
 - move uniformly along this line segment
 - call new location X_n



(diagram from “optGpSampler” paper)

Convergence of Hit-and-Run

- the following is a theorem of Lovász-Vempala from 2004
- let $G \subset \mathbf{R}^d$ be convex, containing a ball of radius r_G , and contained in a ball of radius R_G ; write $\kappa_G := R_G/r_G \geq 1$.
- Then, for some universal $c > 0$, it holds that

$$\gamma_{\text{H\&R}} \geq c \cdot d^{-2} \cdot \kappa_G^{-2}.$$

- high dimension is hard, inhomogeneity of scales is hard

Hit-and-Run Hybrid Slice Sampling

- if π has convex super-level sets, then results of LV give us a bound

$$\sigma_t \geq c \cdot d^{-2} \cdot \kappa_{G(t)}^{-2}$$

- interpretation: life is good if super-level sets $G(t)$ are well-conditioned
 - (if not: worse, though not disastrous)

Well-Conditioned Level Sets

- let $V : \mathbf{R}^d \rightarrow \mathbf{R}$ be m -strongly convex and L -smooth
 - i.e. $\text{eigs}(\text{Hess}V(x)) \in [m, L]$
 - write $\kappa_V = L/m \geq 1$
- let density π have the form $\pi(x) = \text{decreasing}(V(x))$
- then for all t , it holds that $\kappa_{G(t)} \leq \sqrt{\kappa_V}$.

Some Applications

- if π has this form, then $\gamma_H \gtrsim d^{-2} \cdot \kappa_V^{-1} \cdot \gamma_U$
 - H&R-HSS is only worse than ideal SS by factor $d^2 \cdot \kappa_V$
- if e.g. $\pi \propto \exp(-V)$,
 - combine with works on ideal SS, \rightsquigarrow relaxation time of $\lesssim d^3 \cdot \kappa_V$
- if e.g. π is multivariate Student-t, then $\kappa_V = 1$, $\sigma_H \gtrsim d^{-2}$
 - combine with earlier work, \rightsquigarrow relaxation time of $\lesssim d^4$

Some Recap

- slice sampling performs well in theory, and in practice (when possible)
- hybrid slice sampling performs well in practice, is *typically* possible,
 - ... and we provide here some theory to support this
- comparison principles: i) is U good?, ii) is H similar enough to U ?
- generally, $H \preceq U$,
 - ... but if $H_t \succeq \sigma_H \cdot \nu_t$, then $H \succeq \sigma_H \cdot U$.

Some Closing Remarks

- today: exponential rates, Hit-and-Run on the slice
- **in the paper**: slower-than-exponential rates, other examples of on-slice kernels, stepping-out and shrinkage, ‘generalised’ slice sampling with different reference measures,
- theoretical framework is very robust to which on-slice kernels are used
- actually, theoretical framework is **much more general** than slice sampling
 - “Markov chain comparison”, “weak Poincaré inequalities”, ...

— Bonus Material —

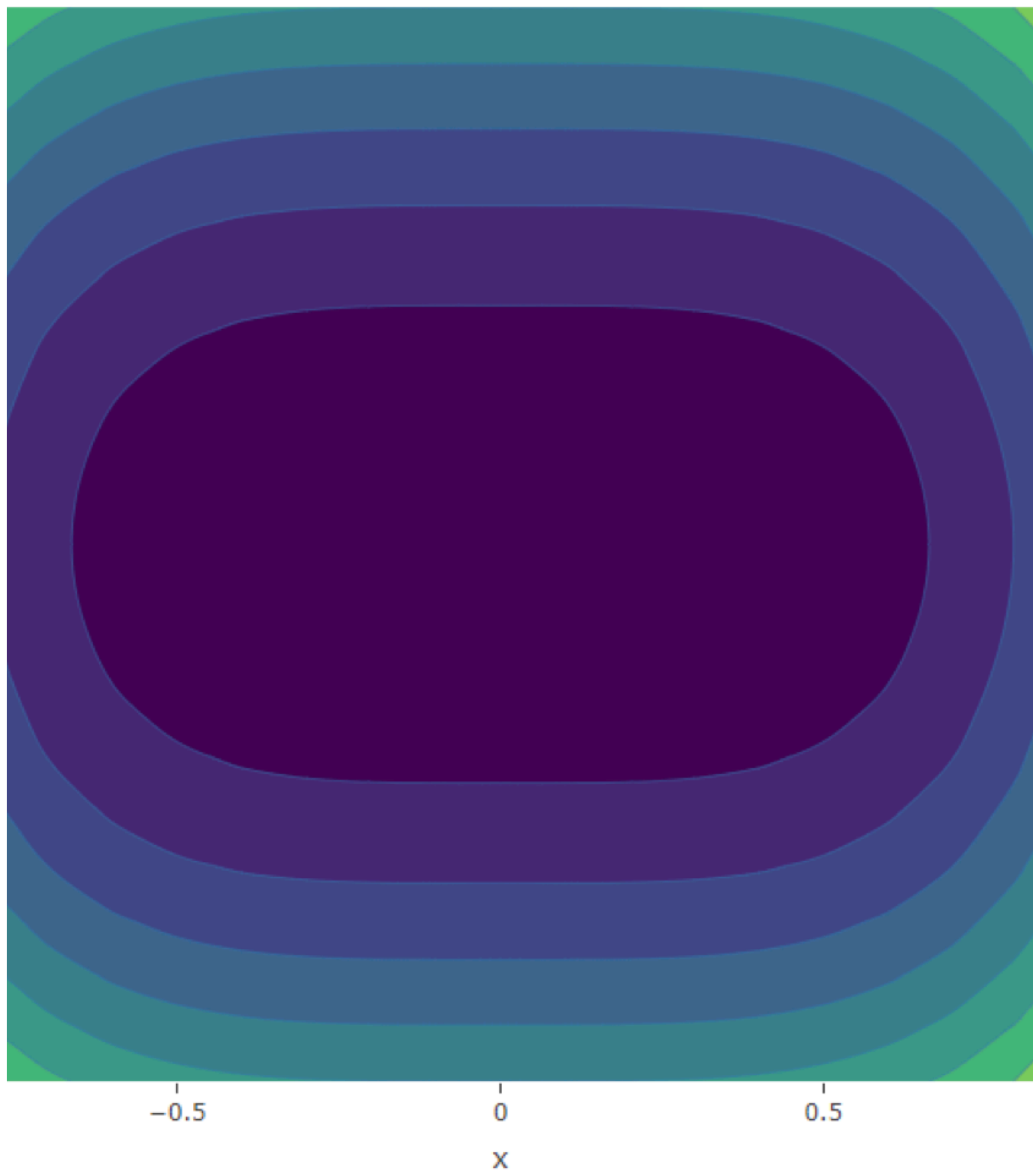
Advanced Applications

- let $1 \leq p_2 \leq p_1$, $1 \leq q_1 \leq q_2$, and suppose that

$$\|x\| \sim 0^+ \implies \|x\|^{p_1} \lesssim V(x) \lesssim \|x\|^{p_2}$$

$$\|x\| \sim \infty \implies \|x\|^{q_1} \lesssim V(x) \lesssim \|x\|^{q_2}$$

- if $p_1 = p_2$, $q_1 \neq q_2$, then convergence rate decays quasi-exponentially
- if $p_1 \neq p_2$, then convergence rate decays only polynomially
- message: in this case, bulk behaviour matters more than tail behaviour



$$\kappa_{\mathbf{G}}(t) \leq \begin{cases} c_{\kappa}^{-} \cdot \left(\log\left(\frac{1}{t}\right)\right)^{\theta} & 0 < t \leq \exp(-1); \\ c_{\kappa}^{+} \cdot \left(\log\left(\frac{1}{t}\right)\right)^{-\vartheta} & \exp(-1) \leq t < 1; \end{cases}$$

with $\theta = \frac{1}{q_1} - \frac{1}{q_2}$, $\vartheta = \frac{1}{p_2} - \frac{1}{p_1}$, and such that the mass function satisfies

$$m(t) \leq c_m \cdot \left(\log\left(\frac{1}{t}\right)\right)^{d/r}$$

with $r = q_1$. By application of Proposition [40](#), we see that for $p_1 = p_2$, there holds a WPI with

$$\beta(s) \leq c^{(1)} \cdot \exp\left(-c^{(2)} \cdot s^{\frac{q_1 \cdot q_1}{q_2 - q_1}}\right),$$

whereas for $p_1 > p_2$, one instead obtains a WPI with

$$\beta(s) \leq c^{(3)} \cdot s^{-\left(1 + \frac{d}{q_1}\right) \cdot \frac{p_1 \cdot p_2}{p_1 - p_2}}.$$

Quantitative Mode Separation

Definition 33. Fix a density function $\varpi : \mathbb{R} \rightarrow \mathbb{T}$, and let $0 < t_1 \leq t_2$ be elements of \mathbb{T} . Say that ϖ is (t_1, t_2) -bimodal if

- for all $t \in \mathbb{T} \setminus [t_1, t_2)$, the super-level set $G(t)$ consists of a single interval, and
- for all $t \in [t_1, t_2)$, the super-level set $G(t)$ consists of a pair of disjoint sub-intervals $G(t) = G_1(t) \sqcup G_2(t)$ such that commonly-labelled sub-intervals are nested, i.e. for $i = 1, 2$ and $t_1 \leq s \leq t < t_2$, there holds the inclusion $G_i(t) \subseteq G_i(s)$.

Moreover, given a (t_1, t_2) -bimodal density ϖ , define $\delta_\varpi : \mathbb{T} \rightarrow [0, \infty)$ by

$$\delta_\varpi(t) = \begin{cases} \text{dist}(G_1(t), G_2(t)) & t \in [t_1, t_2) \\ 0 & \text{otherwise,} \end{cases}$$

where $\text{dist}(G_1(t), G_2(t)) := \inf \{|x - y| : x \in G_1(t), y \in G_2(t)\}$, and write $\Delta_\varpi := \sup \{\delta_\varpi(t) : t \in \mathbb{T}\}$.

Stepping-Out and Shrinkage

Assumption 2. *Let ϖ be a (t_1, t_2) -bimodal density and $h > \Delta_\varpi$ be a stepping-out parameter.*

Under the previous assumption, for $t \in \mathbb{T}$, define the ‘stepping-out and shrinkage’ on-slice kernel with parameter h by

$$H_t(x, A) = \lambda(t) \cdot \nu_t(A) + (1 - \lambda(t)) \cdot \left[\mathbf{1}_{G_1(t)}(x) \cdot \nu_{t,1}(A) + \mathbf{1}_{G_2(t)}(x) \cdot \nu_{t,2}(A) \right],$$

for $x \in \mathbb{R}$, $A \in \mathcal{B}(\mathbb{R})$, where

$$\lambda(t) := \begin{cases} \frac{h - \delta_\varpi(t)}{h} \cdot \frac{m(t)}{m(t) + \delta_\varpi(t)} & t \in [t_1, t_2); \\ 1 & \text{otherwise;} \end{cases}$$
$$\nu_{t,i}(A) := \begin{cases} \frac{\nu(G_i(t) \cap A)}{\nu(G_i(t))} & t \in [t_1, t_2); \\ 0 & \text{otherwise;} \end{cases} \quad \text{for } i = 1, 2.$$

Weak Poincaré Inequalities

Definition 1. We say that a μ -reversible, positive transition kernel P satisfies a weak Poincaré inequality (WPI) if for all $f \in L^2_0(\mu)$ we have

$$\|f\|_\mu^2 \leq s \cdot \mathcal{E}_\mu(P, f) + \beta(s) \cdot \|f\|_{\text{osc}}^2, \quad (3)$$

where $\beta: (0, \infty) \rightarrow [0, \infty)$ is a decreasing function with $\lim_{s \rightarrow \infty} \beta(s) = 0$.

Assumption 1. We assume that for Lebesgue-almost every $t \in \mathbb{T}$, the kernel H_t is ν_t -reversible, positive and satisfies a WPI, i.e. there is a measurable function $\beta: (0, \infty) \times \mathbb{T} \rightarrow [0, \infty)$ with $\beta(\cdot, t)$ satisfying the conditions in Definition 1 for each $t \in \mathbb{T}$, such that for each $s > 0$, $f \in L^2(\nu_t)$,

$$\text{Var}_{\nu_t}(f) \leq s \cdot \mathcal{E}_{\nu_t}(H_t, f) + \beta(s, t) \cdot \|f\|_{\text{osc}}^2. \quad (8)$$

Theorem 11. Under Assumption 1, we have the following comparisons for U and H given in (6) and (7):

For all $f \in L^2(\pi)$,

$$\mathcal{E}(H, f) \leq \mathcal{E}(U, f), \quad (9)$$

and conversely, for all $s > 0$, $f \in L^2(\pi)$,

$$\mathcal{E}(U, f) \leq s \cdot \mathcal{E}(H, f) + \beta(s) \cdot \|f\|_{\text{osc}}^2, \quad (10)$$

where $\beta: (0, \infty) \rightarrow [0, \infty)$ is given by

$$\beta(s) := c^{-1} \cdot \int_{\mathbb{T}} \beta(s, t) \cdot m(t) \, dt.$$

Furthermore, β satisfies the conditions for a WPI in Definition 1.

Metropolis Chains as HSS

4.1 Metropolis chains

Definition 22. Let π a probability measure admitting a density $\varpi = \frac{d\pi}{d\nu}$ with respect to some σ -finite measure ν on G , and let Q be a ν -reversible Markov kernel; we say that such triples (π, ν, Q) are compatible. Define Metropolis (π, ν, Q) to be the Markov kernel P given by

$$P(x, A) = \int_A Q(x, dy) \cdot \alpha(x, y) + \mathbf{1}_A(x) \cdot \bar{\alpha}(x), \quad x \in \mathsf{G}, A \in \mathcal{G},$$

where for $x, y \in \mathsf{G}$,

$$\alpha(x, y) := 1 \wedge \frac{\varpi(y)}{\varpi(x)}, \quad \bar{\alpha}(x) := 1 - \alpha(x), \quad \alpha(x) := \int_{\mathsf{G}} Q(x, dy) \cdot \alpha(x, y).$$

It is known that all such kernels are π -reversible. In what follows, we will use (π, ν, Q) informally to refer to generic compatible triples in the sense described above.

Lemma 27. For any compatible triple (π, ν, Q) , it holds that Metropolis $(\pi, \nu, Q) = \text{HybridSlice}(\pi, \nu, (H_t))$, where

$$H_t(x, A) = Q(x, A \cap \mathsf{G}(t)) + \mathbf{1}_A(x) \cdot Q(x, \mathsf{G}(t)^c).$$

Some Metropolis Chains

Example 23. When $\nu = \text{Leb}$ and Q is a symmetric, ν -reversible kernel, then we can define the Random Walk Metropolis (RWM) kernel,

$$\text{RWM}(\pi, Q) := \text{Metropolis}(\pi, \text{Leb}, Q).$$

It is conventional to work with $Q_\sigma(x, dy) = \mathcal{N}(dy \mid x, \sigma^2 \cdot I_d)$ for some step-size $\sigma > 0$; we will work under this assumption going forward. See also Section 6.3.2 of [29].

Example 24. When ν is a sufficiently-tractable probability measure, we may take $Q(x, \cdot) = \nu$ directly, independently of x . We can thus define the Independent Metropolis–Hastings (IMH) kernel with ‘proposal’ ν ; see [29, Section 6.3.1]:

$$\text{IMH}(\pi, \nu) := \text{Metropolis}(\pi, \nu, \nu).$$

Example 25. When $\nu = \gamma_{\mathbf{m}, \mathbf{C}}$ is a Gaussian measure with mean \mathbf{m} and covariance operator \mathbf{C} , then one may take $\rho, \eta \in (0, 1)$ such that $\rho^2 + \eta^2 = 1$ and define the autoregressive proposal $Q_\eta(x, dy) = \mathcal{N}(dy \mid \mathbf{m} + \rho \cdot (x - \mathbf{m}), \eta \cdot \mathbf{C})$. The resulting Metropolis chain is known as the Preconditioned Crank-Nicolson (pCN) kernel with Gaussian reference $\gamma_{\mathbf{m}, \mathbf{C}}$ and step-size η ; see e.g. [9]:

$$\text{pCN}(\pi, \mathbf{m}, \mathbf{C}, \eta) := \text{Metropolis}(\pi, \gamma_{\mathbf{m}, \mathbf{C}}, Q_\eta).$$

Beyond Hit-and-Run

There are various routes left open by this work; we list here a few.

For one, we have largely focused on Simple Slice Sampling as the ideal algorithm, and using Hit-and-Run as on-slice kernels, due to their relative genericity and ease of implementation. In specific settings, other on-slice kernels are of substantial interest. For super-level sets with some coordinate-related structure, Gibbs sampling (also ‘Coordinate Hit-and-Run’) is a natural option, for which results have recently been obtained in the convex setting; see [17]. For super-level sets which take the form of polytopes, yet more on-Slice Samplers are available, including the Dikin walk [16], the Vaidya and John walks [7], and various gradient-based samplers which introduce additional geometric structure to the problem (e.g. [20]); some of these methods come with theoretical guarantees in the form of estimates on the conductance or spectral gap of the kernel, which can readily be used in our framework.