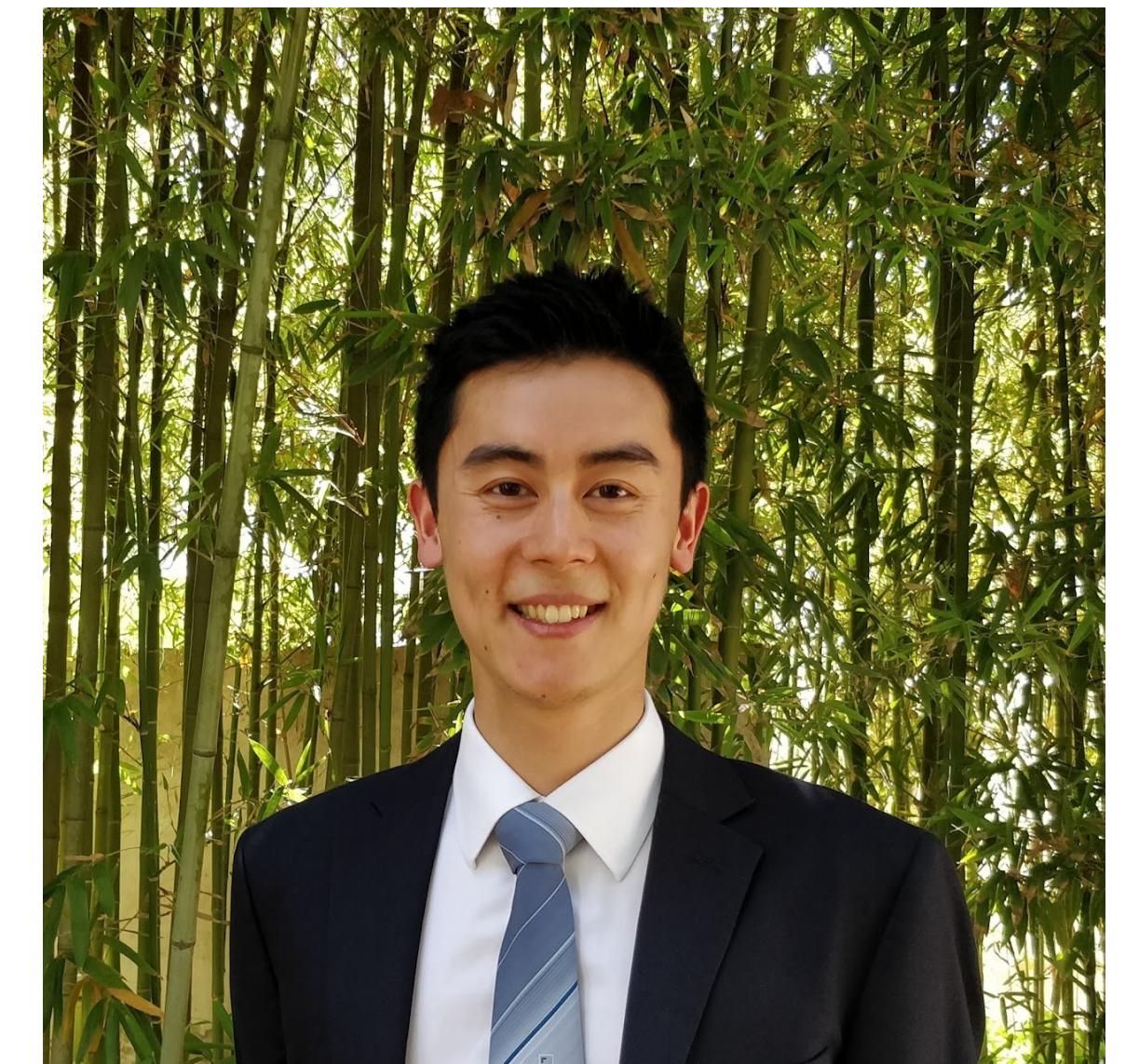


# Comparison Theorems for Slice Sampling

Probability Seminar, University of Bristol, 15 November 2024

Sam Power, University of Bristol

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



# Some personal background

- Sam Power @ University of Bristol
  - Lecturer in Statistical Science here since January 2024
  - Postdoc with Christophe Andrieu, Anthony Lee here since August 2020
- Interested in { Modelling, Computation, Analysis, Probability, ... }
- often Markov Processes (Long-Time Behaviour, Comparisons)

# Plan for today

- a sampling problem
- a cute reduction / trick
- a Markov chain
- a convergence analysis (not ours)
- a refined analysis (ours)
- some illustrative examples

**feel free to stop me at any point**

# An Elementary Trick

- let  $\pi \ll \text{Leb}$ ; write  $\pi(dx) = \pi(x) dx$
- observe that

$$\pi(x) = \int_{\mathbf{R}_+} \mathbf{1} [\pi(x) > t] dt$$

- introduce

$$\Pi(dx, dt) = \mathbf{1} [\pi(x) > t] dx dt$$

- if you ‘solve’  $\Pi$ , then you ‘solve’  $\pi$

# A Trade

- if you ‘solve’  $\Pi$ , then you ‘solve’  $\pi$
- there ought to be *some* ‘conservation of difficulty’, at some level
- general density on a simple set  $\rightsquigarrow$  uniform density on a general set
- at least *some* tasks ought to be easier in this world

# ‘Vertical Likelihood’

- in some circles, this is known as the ‘vertical likelihood representation’
- one can build various methods on the back of this representation
  - { Accept / Reject, Nested Sampling, Wang-Landau, MUCA, ...}
- today: ‘Markov chain Monte Carlo’ on  $\Pi$

# Markov chain Monte Carlo

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

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

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

# Coming up with a Markov chain

- $\Pi(dx, dt) = 1[\pi(x) > t] dx dt$
- $(x, t)$  has a natural bipartite structure
- simple option: alternating conditional simulation

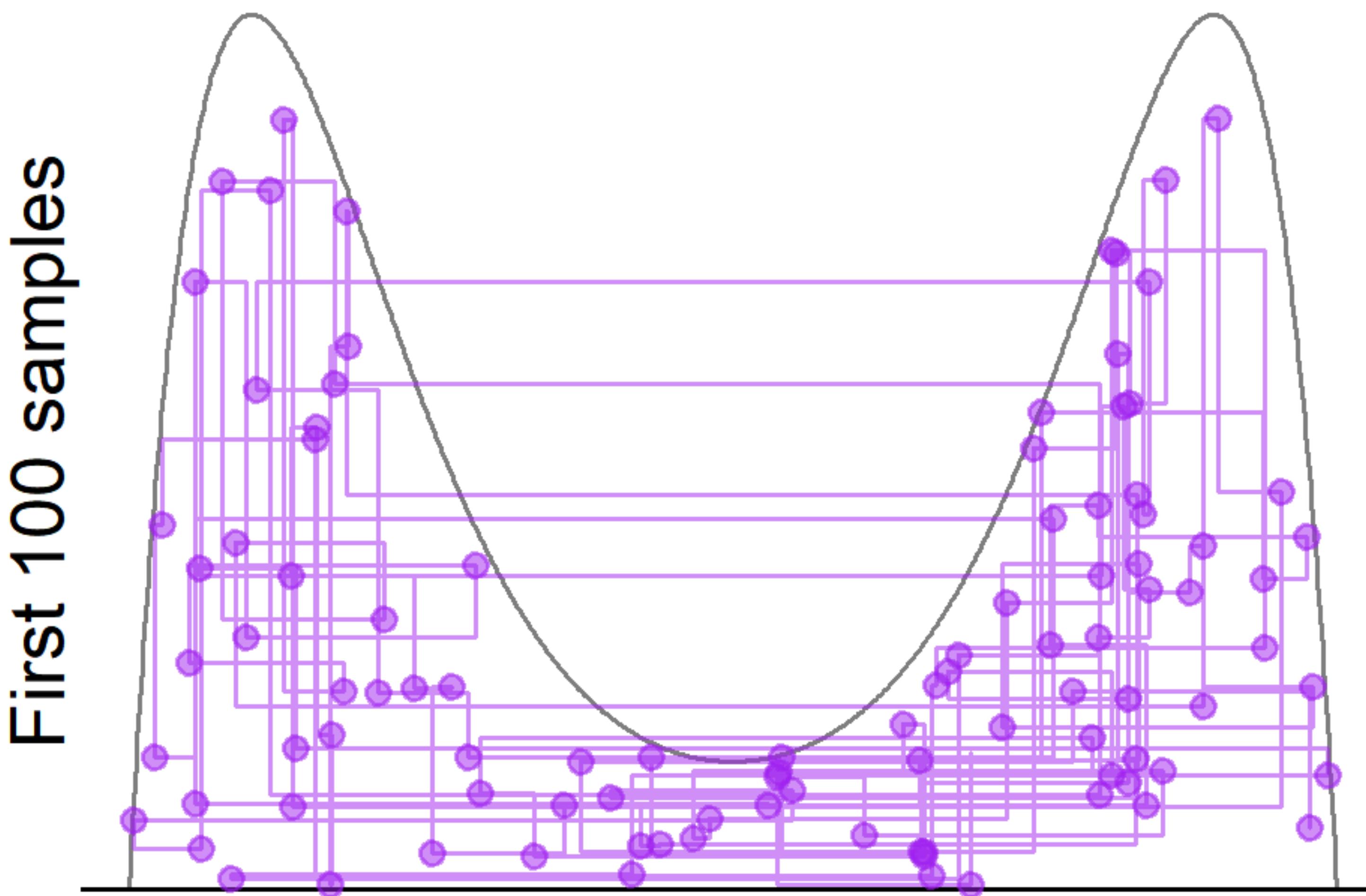
$$[T_n \mid X_{n-1} = x] \sim \Pi(dt \mid X = x)$$

$$[X_n \mid T_n = t] \sim \Pi(dx \mid T = t)$$

- ‘two-component Gibbs sampler’

# What do these conditionals look like?

- $\Pi(dt | X = x)$  is just  $\text{Unif}\left(dt; [0, \pi(x)]\right)$
- $\Pi(dx | T = t)$  is just  $\text{Unif}\left(dx; G(t)\right) \dots$ 
  - ... where  $G(t) = \{x : \pi(x) > t\}$
  - also, will need  $m(t) = \text{Vol}(G(t))$ ,  $\nu_t = \text{Unif}\left(dx; G(t)\right)$
  - (time to get the chalk out!)



# Neal's ‘Slice Sampler’

# Some Connections

- the use of the ‘vertical likelihood’ representation to construct a Markov process has antecedents
  - ‘data augmentation’, imputation strategies for missing data problems
  - ‘random cluster model’ for Ising / Percolation
    - { Fortuin-Kasteleyn / Swendsen-Wang / Edwards-Sokal / ... }
    - other { dual / loop-type / ... } representations of spin systems

# Long-Time Behaviour of Slice Sampling

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

# Detailed Study of Dynamics

- for accurate quantitative understanding, need to study the chain carefully
  - we care about the chain  $(X_n)$ ,
  - but we implement the chain  $(X_n, T_n)$ ,
  - and we might even prefer to look at  $(T_n)$ , which is univariate
- duality / interleaving: all three processes converge equivalently

# Features of the Height Process

- the chain  $(T_n)$  depends on  $\pi$  only through the mass function  $m(t)$ :

$$P(t, t') = \frac{m(t'')}{t''} - \int_{t''}^{\infty} \frac{m(s)}{s^2} ds, \quad t'' = \max \{t, t'\}$$

- so, we can always study a ‘nicer’  $\tilde{\pi}$  for which

$$m_\pi = m_{\tilde{\pi}}$$

- WLOG { centered, isotropic, spherically-symmetric, unimodal, ... }

# The Art of WLOG

- fix a  $\pi$  (hence  $m$ ), pick a dimension  $d \in \mathbf{N}$ , and set

$$\phi(r) = \log\left(\frac{1}{m^{-1}(rd)}\right),$$

which is increasing in  $r$

- then, with  $\tilde{\pi}(x) \propto \exp(-\phi(\|x\|))$ , we have  $m_\pi = m_{\tilde{\pi}}$ 
  - $\approx$  symmetric rearrangement; very regular

# Convergence by Coupling

- given positions  $X, X'$ , sample heights by  $U \sim \text{Unif}([0,1])$ , then

$$T = U \cdot \pi(X), \quad T' = U \cdot \pi(X')$$

- match relative heights
- given heights  $T, T'$ , sample positions by  $S \sim \text{Unif}(\mathbf{S}^{d-1})$ , then

$$X = \phi^{-1} \left( \log \left( \frac{1}{T} \right) \right)^{1/d} \cdot S, \quad X' = \phi^{-1} \left( \log \left( \frac{1}{T'} \right) \right)^{1/d} \cdot S$$

- match direction and magnitude

# Convergence Rates

- for convex  $\phi$ , chain contracts at rate  $\leq \frac{d}{d+1}$
- for  $\phi(r) = \frac{d+m}{2} \cdot \log(1+r^2)$ , chain contracts at rate  $\leq \frac{d \cdot (d+m)}{(d+1) \cdot (d+m-1)}$
- other specific cases can be analysed
- byproduct of coupling construction: chain is *stochastically monotone*

# Some General Comments

- increasing dimension slows down convergence
- heavy tails slow down convergence a bit
  - ... but not *too* dramatically
- very hard to make the convergence awful
  - by “WLOG”, always equivalent to a ‘pretty reasonable’  $\pi$
- impression: method seems robust, potentially algorithmically appealing

# Partial Recap

- sampling problem
- vertical likelihood representation
- Slice Sampling
- practical algorithm
  - ... ?

# Implementing Slice Sampling

$$\text{“}\left[ X_n \mid T_n = t \right] \sim \Pi \left( dx \mid T = t \right)\text{”}$$

– Sam Power, Slide 11

# Life on the Slice

- recall that  $\Pi(dx \mid T = t)$  is  $\nu_t = \text{Unif}(dx; G(t))$
- if  $G(t)$  is a { ball, box, simplex, ... }, then sampling from  $\nu_t$  is fine
- if not, then we have a new problem

# Hybrid Slice Sampling

- instead of
  - “ $[X_n \mid T_n = t] \sim \nu_t$ ”
- select a  $\nu_t$ -invariant Markov kernel  $H_t$ , and do
  - “ $[X_n \mid T_n = t, X_{n-1} = x] \sim H_t(x, dx')$ ”

# Towards Comparison

- let  $U(x, dx')$  denote the “Ideal” Slice Sampling kernel
- let  $H(x, dx')$  denote the “Hybrid” Slice Sampling kernel
- how do these kernels *compare* in terms of efficient convergence?

# Dirichlet Form Analysis of Markov Chains

- if  $P$  is a  $\mu$ -reversible Markov kernel, then define

$$\mathcal{E}_P(f) = \int \mu(dx) \cdot P(x, dy) \cdot \frac{[f(x) - f(y)]^2}{2}$$

- note that taking  $P = \mu$  gives  $\mathcal{E}_\mu(f) = \text{var}_\mu(f)$
- if  $\mathcal{E}_P(f) \geq \gamma \cdot \mathcal{E}_\mu(f)$  for all  $f$ , then “spectral gap”  $\geq \gamma$ 
  - “energy  $\approx$  entropy dissipation  $\gtrsim$  entropy”
  - weakened / generalised version:  $\mathcal{E}_\mu(f) \leq \alpha \cdot \mathcal{E}_P(f) + \beta \cdot \Phi(f)$

# Comparisons for HSS

- in complete generality,  $\mathcal{E}_H \leq \mathcal{E}_U$ 
  - the ideal Slice Sampler should be preferred when possible
- we want to examine when  $\mathcal{E}_H \gtrsim \mathcal{E}_U$ 
  - the practical algorithm is almost as good as the theoretical ideal

*when should this be true?*

# A Natural Approach

- 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)

# Approximate ... how?

- HSS replaces “ $X \sim \nu_t$ ” by “ $X \sim H_t(x, dx')$ ”
- to study how  $H$  approximates  $U$  ...
  - ... first study how  $H_t$  approximates  $\nu_t$
- assumption: for each  $t$ , there exists  $\sigma_t > 0$  such that  $\mathcal{E}_{H_t} \geq \sigma_t \cdot \mathcal{E}_{\nu_t}$
- “on-slice kernel” mixes to “uniform on-slice” at rate  $\sigma_t$

# A Generic Result

- suppose that
  - for each  $t$ ,  $\mathcal{E}_{H_t} \geq \sigma_t \cdot \mathcal{E}_{\nu_t}$ , and
  - $\sigma_H := \inf_{t \in T} \sigma_t > 0$
- it then follows that  $\mathcal{E}_H \geq \sigma_H \cdot \mathcal{E}_U$
- moreover, if  $\mathcal{E}_U \geq \gamma_U \cdot \mathcal{E}_\pi$ , then it further follows that  $\mathcal{E}_H \geq \sigma_H \cdot \gamma_H \cdot \mathcal{E}_\pi$
- interpretation: HSS is at most a factor  $\sigma_H$  ‘worse’ than ideal SS

# Some Remarks

- this particular analysis is quite worst-case, but easy to state
- our methods also provide more refined estimates
  - can accommodate inhomogeneous mixing rates
  - can accommodate on-slice kernels with slower-than-exponential mixing

# A Key Decomposition

- key decomposition formula:

$$\mathcal{E}_H(f) = \int \mathcal{E}_{H_t}(f) \cdot m(t) dt$$

- taking  $H_t = \nu_t$  yields

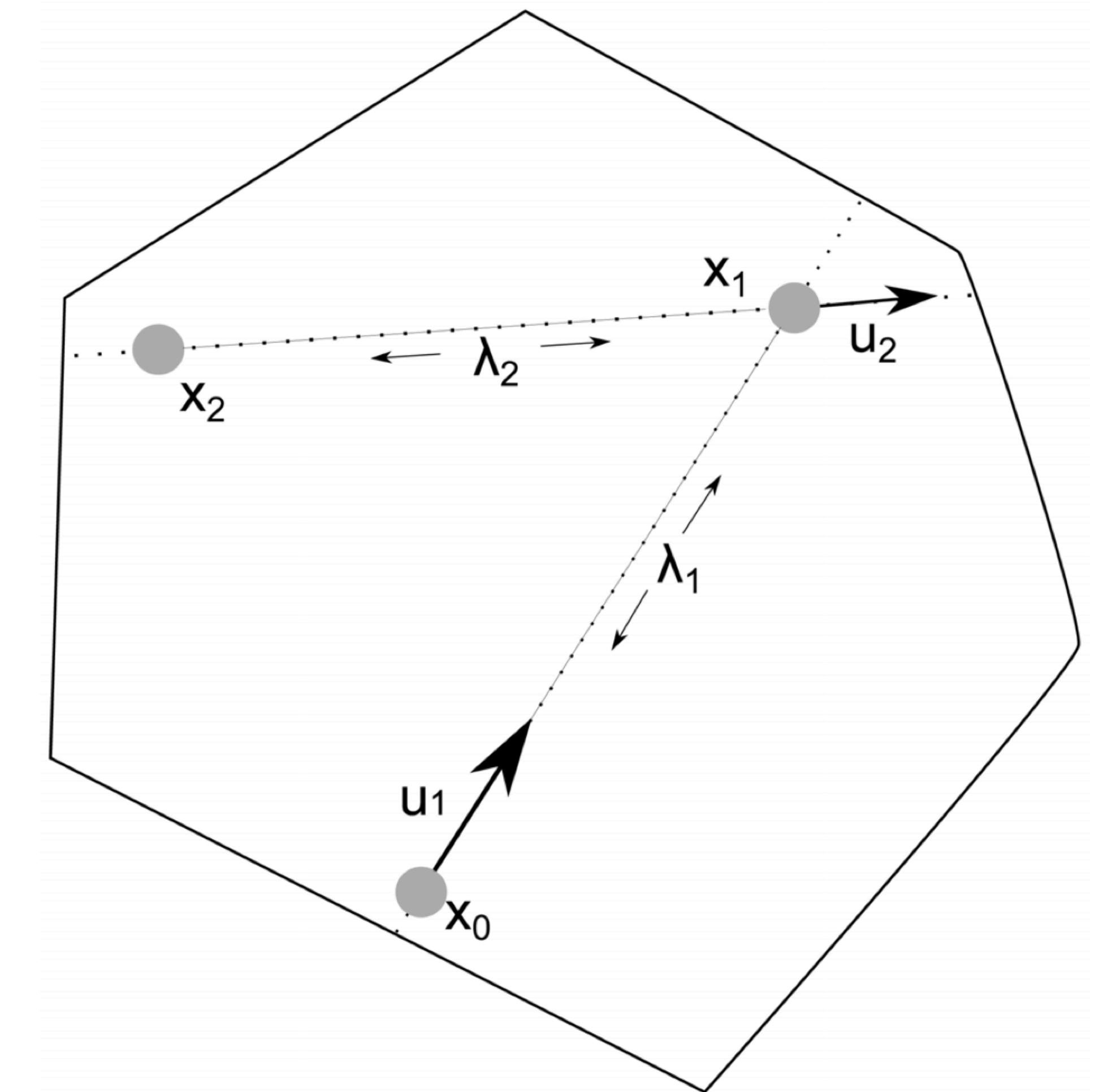
$$\mathcal{E}_U(f) = \int \mathcal{E}_{\nu_t}(f) \cdot m(t) dt$$

- easy to combine with estimates like

$$\mathcal{E}_{\nu_t}(f) \leq \alpha \cdot \mathcal{E}_{H_t}(f) + \beta \cdot \Phi(f)$$

# Case Study: Hit-and-Run on the Slice

- low-tech 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  $L_n = (x + U_n \mathbf{R}) \cap G$
  - sample  $X_n \sim \text{Unif}(L_n)$
- easy to implement, well-studied, ...



(diagram from “optGpSampler” paper)

# Convergence of Hit-and-Run

- let  $G \subset \mathbf{R}^d$  be convex,
  - containing a ball of radius  $r_G$ , contained in a ball of radius  $R_G$  (chalk!)
  - write  $\kappa_G := R_G/r_G \geq 1$ .
- Lovász-Vempala prove that 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  $\pi$  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  $\pi$  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  $\pi$  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,  $\rightsquigarrow$  decorrelation time of  $\lesssim d^3 \cdot \kappa_V$
- if e.g.  $\pi$  is multivariate Student-t, then  $\kappa_V = 1$ ,  $\sigma_H \gtrsim d^{-2}$ 
  - combine with earlier work,  $\rightsquigarrow$  decorrelation 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  $H \succeq \sigma_H \cdot U$ .

# A Glimpse at the Paper

- 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”, ...

# Finale

- Vertical Likelihood Representation  $\rightsquigarrow$  Slice Sampling Markov chain
- Practical considerations  $\rightsquigarrow$  *Hybrid* Slice Sampling Algorithm
- Dirichlet Forms  $\rightsquigarrow$  Comparison Theory for Markov Chains
- Mixing of On-Slice Kernels  $\rightsquigarrow$  Complexity Bounds for HSS



# **– Bonus Material –**

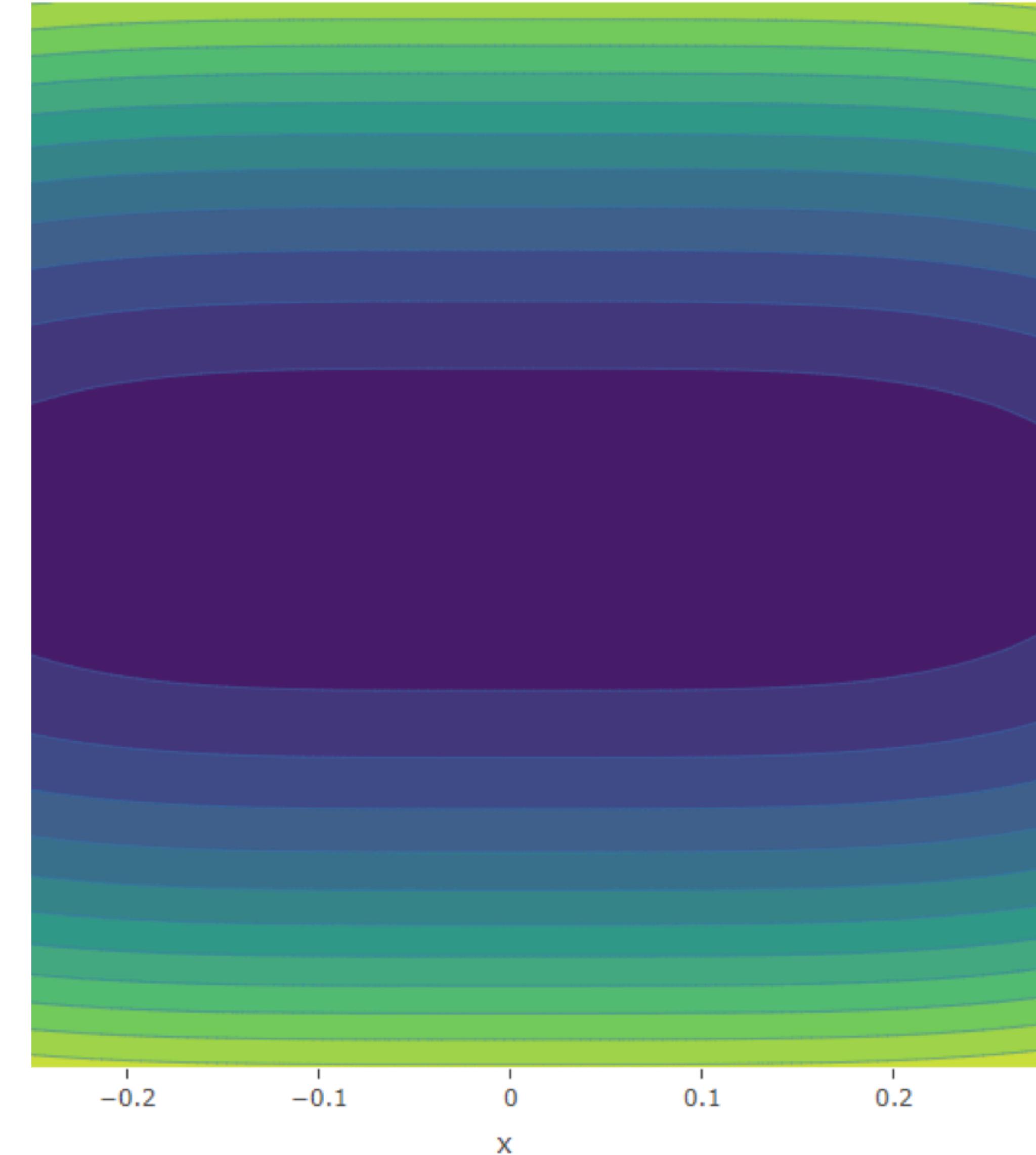
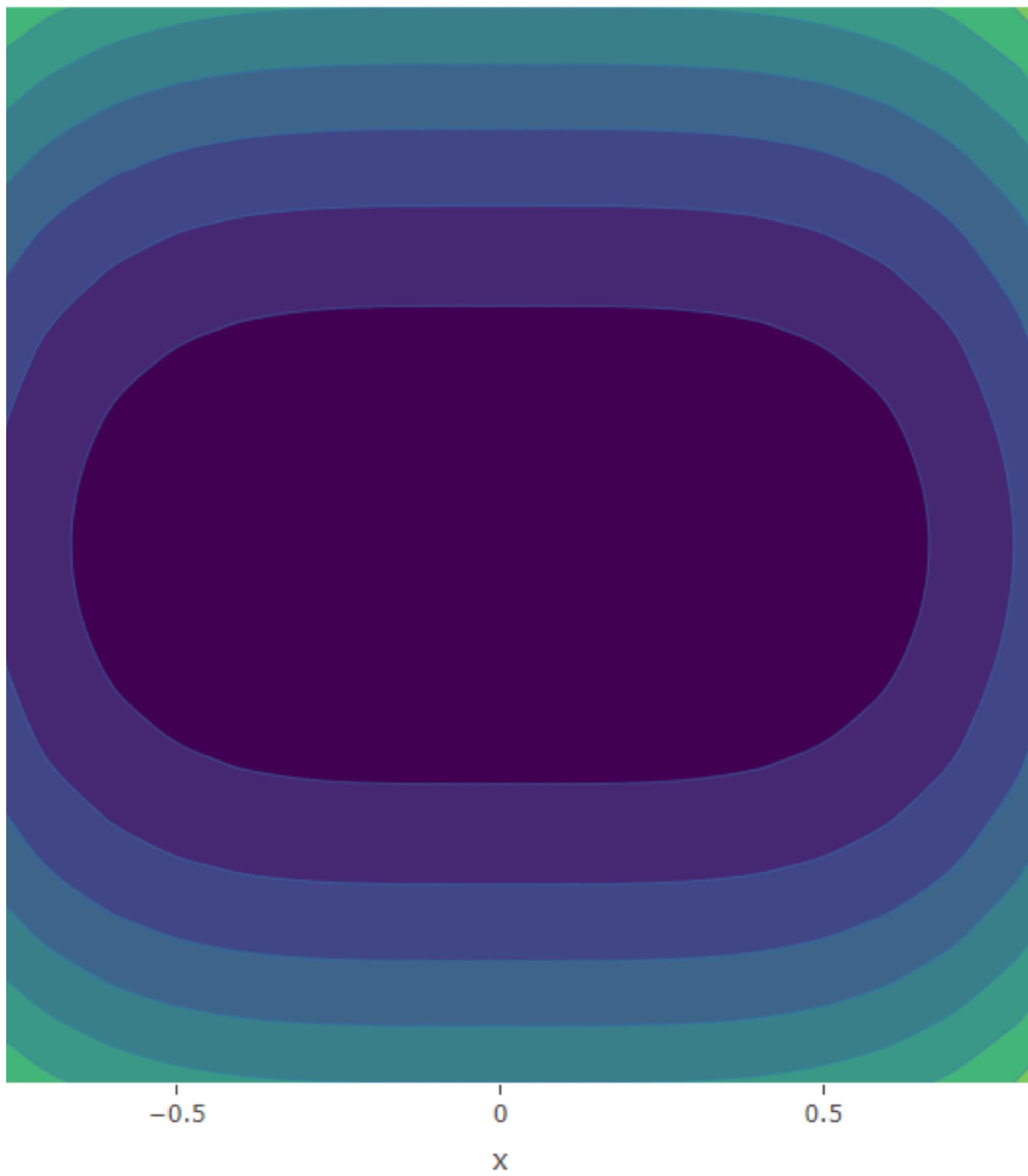
# Advanced Applications

- 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



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

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  $\varpi : \mathbb{R} \rightarrow \mathsf{T}$ , and let  $0 < t_1 \leq t_2$  be elements of  $\mathsf{T}$ . Say that  $\varpi$  is  $(t_1, t_2)$ -bimodal if

- for all  $t \in \mathsf{T} \setminus [t_1, t_2]$ , the super-level set  $\mathbf{G}(t)$  consists of a single interval, and
- for all  $t \in [t_1, t_2]$ , the super-level set  $\mathbf{G}(t)$  consists of a pair of disjoint sub-intervals  $\mathbf{G}(t) = \mathbf{G}_1(t) \sqcup \mathbf{G}_2(t)$  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  $\mathbf{G}_i(t) \subseteq \mathbf{G}_i(s)$ .

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

$$\delta_\varpi(t) = \begin{cases} \text{dist}(\mathbf{G}_1(t), \mathbf{G}_2(t)) & t \in [t_1, t_2] \\ 0 & \text{otherwise,} \end{cases}$$

where  $\text{dist}(\mathbf{G}_1(t), \mathbf{G}_2(t)) := \inf \{|x - y| : x \in \mathbf{G}_1(t), y \in \mathbf{G}_2(t)\}$ , and write  $\Delta_\varpi := \sup \{\delta_\varpi(t) : t \in \mathsf{T}\}$ .

# Stepping-Out and Shrinkage

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

Under the previous assumption, for  $t \in T$ , 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

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

# Weak Poincaré Inequalities

**Definition 1.** We say that a  $\mu$ -reversible, positive transition kernel  $P$  satisfies a weak Poincaré inequality (WPI) if for all  $f \in L_0^2(\mu)$  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  $t \in T$ , the kernel  $H_t$  is  $\nu_t$ -reversible, positive and satisfies a WPI, i.e. there is a measurable function  $\beta: (0, \infty) \times T \rightarrow [0, \infty)$  with  $\beta(\cdot, t)$  satisfying the conditions in Definition 1 for each  $t \in T$ , 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

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

Furthermore,  $\beta$  satisfies the conditions for a WPI in Definition 1.

# Metropolis Chains as HSS

## 4.1 Metropolis chains

**Definition 22.** Let  $\pi$  a probability measure admitting a density  $\varpi = \frac{d\pi}{d\nu}$  with respect to some  $\sigma$ -finite measure  $\nu$  on  $\mathsf{G}$ , and let  $Q$  be a  $\nu$ -reversible Markov kernel; we say that such triples  $(\pi, \nu, Q)$  are compatible. Define  $\text{Metropolis}(\pi, \nu, Q)$  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  $\pi$ -reversible. In what follows, we will use  $(\pi, \nu, Q)$  informally to refer to generic compatible triples in the sense described above.

**Lemma 27.** For any compatible triple  $(\pi, \nu, Q)$ , it holds that  $\text{Metropolis}(\pi, \nu, Q) = \text{HybridSlice}(\pi, \nu, (H_t))$ , 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,  $\nu$ -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  $Q_\sigma(x, dy) = \mathcal{N}(dy | x, \sigma^2 \cdot I_d)$  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  $\nu$  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’  $\nu$ ; see [29, Section 6.3.1]:

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

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

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

# 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.