# Comparison Theorems for Practical Slice Sampling

Slice Sampling is a popular algorithm for approximate sampling from intractable probability ~~distributions~~ distributions, used in {JAGS, Matlab, ...}.

Its popularity stems from its wide applicability, and robustness, in both practice and theory. [intuitive geometric formulation]

An outstanding theoretical challenge has been that while the ideal slice sampler ~~has excellent~~ admits a robust and elegant convergence theory, practical implementations ~~████~~ typically involve additional approximations, ~~████~~ which prevents the existing theory from holding as is.

In recent work, we develop a theoretical framework for the analysis of such "hybrid" slice samplers, ~~████~~ facilitating novel convergence results for slice sampling as implemented in practice.

We provide a number of concrete examples which illustrate the flexibility and practicality of our approach, including i) stepping-out and shrinkage procedures, ii) hit-and-run "on the slice",

No prior knowledge of the slice sampling algorithm will be assumed, and relevant notions of Markov chain convergence ~~will be~~ will be ~~developed~~ introduced as appropriate in the talk.

## Optimality of MLE

in what sense? $\longrightarrow$ MLE is typically consistent, asymptotically unbiased, good variance.

want a finite-$n$ comparison.

simplify: think about unbiased estimators; what is possible?
(good estimators might be close to unbiased)
( c.f. Gauss-Markov).
mathematically, convenient family.

Start with
$$\begin{cases} \int P_\theta(x)\, dx = 1 \\ \int P_\theta(x)\, \hat{\theta}(x)\, dx = \theta \end{cases} \Rightarrow \theta\text{-independent relation.}$$

general $\underline{\hspace{2cm}} \int P_\theta(x)\, F(x) = G(\theta).$

define $F(\theta) = \int P_\theta(x)\, f(x)\, dx \in \mathbb{R}^Q$

$$\frac{\partial F}{\partial \theta_k} = \frac{\partial}{\partial \theta_k} \int P_\theta(x)\, f(x)\, dx$$

$$= \int P_\theta(x)\, \frac{\partial}{\partial \theta_k} \log P_\theta(x) \cdot f(x)\, dx$$

$$\nabla_\theta F(\theta) = \text{Cov}_\theta\left(\nabla_\theta \log P_\theta(x), f(x)\right) \in \mathbb{R}^P.$$

$f(x) = 1 \Rightarrow 0 = \text{Cov}_\theta \nabla$

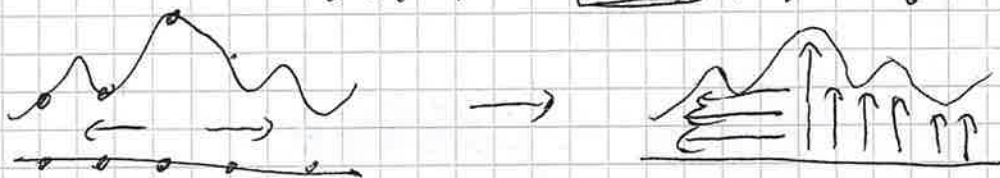① **MCMC** $\qquad \pi \to \{x_1 \to x_2 \to \cdots$

$$\mathrm{Law}(x_t) \to \pi$$

$$\frac{1}{T}\sum_{t=1}^{T} f(x_t) \to \int \pi(dx) f(x).$$

oracle $\log \pi$ (up to a constant), [more]

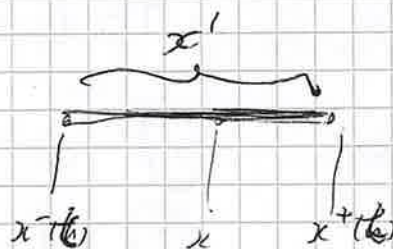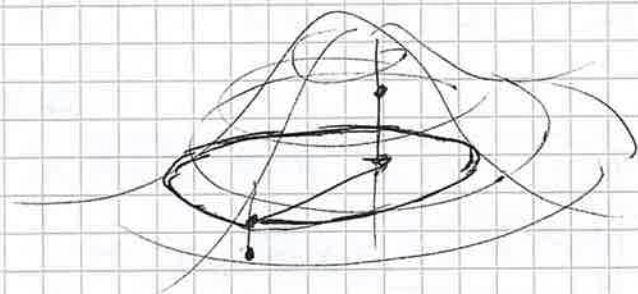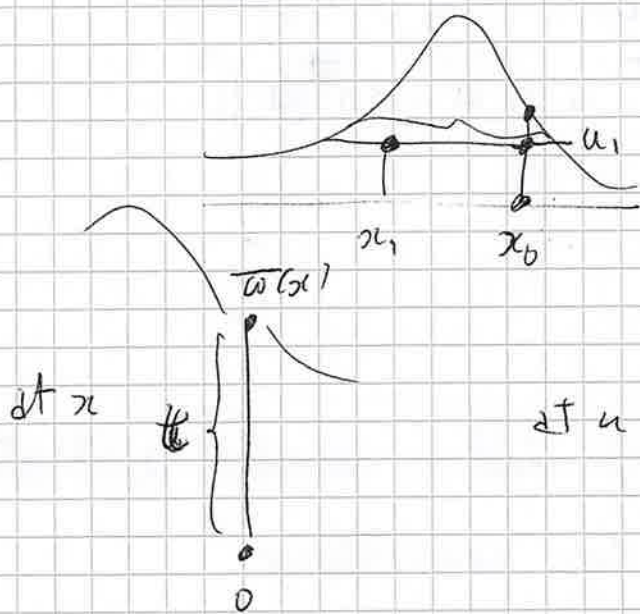② **Slice Sampling** $\qquad \pi(dx) = \varpi(x)\, \lambda^d(dx)$

$$\Pi(d(x, t)) = \boxed{\phantom{xxx}} \, \mathbb{I}[\,0 \le t \le \varpi(x)\,]\, \lambda^{d+1}(d(x,t))$$



"iterated conditional simulation":

① at $x$, $\quad u \sim \Pi(dt|x) = \mathrm{Unif}([0, \varpi(x)])$

(sample a height)

② at $t$, $\quad x \sim \Pi(dx|t) = \mathrm{Unif}(\{x : \varpi(x) \ge t\})$.

(sample a point of the height)



$$x \to \{ \longrightarrow \cdots \quad x'$$

$$x' \sim U(x, dx').\qquad [\text{Markov kernel}]$$

①

$$\left(\begin{smallmatrix} a \\ b \end{smallmatrix}\right)^T - \left(\begin{smallmatrix} a \\ b \end{smallmatrix}\right) = a^T I(\theta) a + 2 a^T b + b^T \text{cov}(\hat{\theta}(x)) b$$

$$\frac{d}{da}: \quad I(\theta) a + b = 0, \quad a = -I(\theta)^{-1} b.$$

$$\Rightarrow \quad b^T \left(\text{cov } \hat{\theta}(x) - I(\theta)\right) b \geq 0.$$

$$\left( -I(\theta)^{-1} \text{\reflectbox{V}} \right)$$

Caveats $\qquad$ let $\qquad \mathbb{E}\,\hat{\theta}(x) = (\cancel{A} (1-\varepsilon)) \theta$

$$\mathbb{E} \| (1-\varepsilon) \hat{\theta}(x) - \theta \|^2 = (1-\varepsilon)^2 \mathbb{E} \| \hat{\theta} \|^2$$

$$\cancel{\text{cov}\,(1-\varepsilon)\theta}$$

$$\text{cov}\left(A\hat{\theta}(x)\right) = A \,\text{cov}(\hat{\theta}(x)) A^T$$

$$\mathbb{E}\left(A\hat{\theta}(x) - \theta\right)^{\otimes 2} = \mathbb{E}\left(A(\hat{\theta}(x) - \theta)\right) + (A\theta)^{\otimes 2}$$

$$= A C A^T + \|(A-1)\theta\|^2$$

$$\boxed{\cancel{\frac{\partial}{\partial A} \quad A \to A + \delta A \quad (\delta A) C A^T}}$$

Bias–Variance Tradeoff

One approach $\qquad$ ① come up with + sensible estimator

$\qquad\qquad\qquad$ ② check for consistency

$\qquad\qquad\qquad$ ③ check that bias is sub-dominant

$\qquad\qquad\qquad$ ④ apply shrinkage/regularisation
$\qquad\qquad\qquad\qquad$ to control variance

$\qquad\qquad$ (other stuff)

$\text{Thm}^{(N)}$ SS leaves $\pi$ invariant $\implies$ valid for MCMC

$\text{Thm}^{(RR)}$ weak conditions $\Rightarrow$ SS is geometrically ergodic $\Rightarrow$ more thms

observation   I know no example where SS is <u>not</u> geometrically ergodic

Thm (NES).  $\pi$ spherically symmetric, log-concave
$$\implies \text{convergence is like } e^{-t/d}, \ t \text{ iterations}$$

also: $\pi$ itself doesn't matter, just mass of level sets

Thm (Sch)  $\pi$ = student-$t$ $\implies$ convergence like $e^{-t/d^2}$

(heavy tails but still geometric)

$\exists$ other examples.

$\implies$ When implementable, SS is robust and
handles dimensionality well, for a $0^{th}$-order method.

<u>Issue</u>   "sample $x \sim \text{Unif}(\{x: \sigma(x) > t\})$
$$= \text{Unif}(G(t))$$

· If $G(t) \in \bigcirc, \square, \cdots$, sure
· If $G(t)$ is more arbitrary .... work for it.

Let $\nu_t(dx) = \text{Unif}(dx; G(t))$

SS:   $\begin{aligned} t &\sim \text{Unif}([0, \sigma(x)]) \\ x' &\sim \nu_t \end{aligned}$

HSS: replace $x' \sim \nu_t$
with   $x' \sim \text{MCMC}(x \to x'; \text{target} = \nu_t)$

Usually easier to implement, slower Markov chain.
(examples!)

<u>Question</u> :  how much slower?
how do we trade off ease/efficiency in practice?
what price are we currently paying?

"comparison theory"

Fixed $p$, $n \to \infty$: bias not always an issue.
$p \asymp n$, $(p,n) \to (\infty, \infty)$: more subtle.

$$MSE(\hat{\theta}(x)) = \| \mathbb{E}\,\hat{\theta}(x) - \theta \|_2^2 + Tr(Cov(\hat{\theta}(x)))$$

statistical "efficiency" : $Var \approx Var_{CRLB}$
(asymptotic)

(often only asymptotically unbiased)
(in practice, want "bias is not an issue")

## Proof Elements

① $\int P_\theta(x) \nabla_\theta \log P_\theta(x)\, dx = 0.$

② $\int P_\theta(x) \nabla_\theta^2 -\log P_\theta(x)\, dx$
$$= \int P_\theta(x) \left(\nabla_\theta \log P_\theta(x)\right)\left(\nabla_\theta \log P_\theta(x)\right)^\top dx$$

③ $\int P_\theta(x) \log \dfrac{P_{\theta_*}(x)}{P_\theta(x)}\, dx \geq 0.$

④ $\ell(\theta, X) = \log P_\theta(X)$
$\ell(\theta) = \mathbb{E}_{\theta_*}[\log P_\theta(x)] \leq \ell(\theta_*).$

## Proof

let $\hat{\theta}(x)$ be unbiased for $\theta$.
$$\Rightarrow \forall \theta \in \Theta, \quad \int P_\theta(x)\,\hat{\theta}(x)\, dx = \theta.$$

$\dfrac{\partial}{\partial \theta_i} \Rightarrow \int P_\theta(x) \left(\dfrac{\partial}{\partial \theta_i} \log P_\theta(x)\right) \hat{\theta}_j(x)\, dx = \delta_{ij}$

$$Cov(\nabla_\theta \log P_\theta(x), \hat{\theta}(x)) = I.$$

$$Cov \begin{pmatrix} \nabla_\theta \log P_\theta(x) \\ \hat{\theta}(x) \end{pmatrix} = \begin{pmatrix} \mathcal{I}(\theta) & I_p \\ I_p & Cov(\hat{\theta}(x)) \end{pmatrix} \succeq 0.$$

$\begin{pmatrix} v_1 \\ v_2 \end{pmatrix}^\top$ ... $\begin{pmatrix} v \\ v \end{pmatrix}^\top \begin{pmatrix} 0 & v \\ v & Cov(\hat{\theta}(x)) v \end{pmatrix} =$

## Convergence of Markov chains

$$\left( \begin{array}{l} \mu_n = \mu P^n. \\[2mm] d(\mu, \pi) = \int \pi(dx) \left( \frac{d\mu}{d\pi}(x) - 1 \right)^2 \\[4mm] \qquad = \left\| \frac{d\mu}{d\pi} - 1 \right\|^2_{L^2(\pi)}. \end{array} \right)$$

$$\| f \|^2_{L^2(\pi)} = \int \pi(dx) \, f(x)^2$$

$$Pf(x) = \int P(x, dy) f(y).$$

$$P^n f(x) \longrightarrow \pi(f) \qquad (\text{ergodicity})$$

$$\| P^n f(x) - \pi(f) \|^2_{L^2(\pi)} \longrightarrow 0 \qquad \text{"}L^2 \text{ convergence"}$$
$$(\text{focus on } \pi(f) = 0)$$

$\simeq$ <u>Best case</u> : $\quad \| P f \|^2_2 \leq (1 - \lambda) \| f \|^2_2 \qquad \forall f \in L^2_0(\pi)$

$$\Rightarrow \text{ exponential rate.}$$

actually, equivalent to

$$\mathcal{E}(P^* P, f) \geqslant \lambda \| f \|^2_2$$

$$\frac{1}{2} \int \pi(dx) \overbrace{(P^* P)(x, dy)} (f(x) - f(y))^2 \quad \Big\uparrow \begin{array}{l}\text{remaining} \\ \text{energy}\end{array}$$

$$\underbrace{(I - P^2) = (I + P)(I - P) \geq (I - P)}_{\substack{\text{energy} \\ \text{dissipation}}}$$

<u>Fact</u> $P$ "positive", $\mathcal{E}(P, f) \geqslant \gamma \| f \|^2$
     reversible
$$\Rightarrow \quad \| P^n f \|^2 \leq (1 - \gamma)^k \| f \|$$

"good energy dissipation $\Rightarrow$ good convergence" ③

# CRAMÉR-RAO LOWER BOUND

Gauss-Markov : Estimator $\wedge$ Linear $\wedge$ Unbiased
$$\Rightarrow \text{OLS is minimum variance}$$

Same for MLE?

$$\ell_n(\theta; x_{1:n}) = \log P_\theta(x_{1:n})$$

$$\widehat{\theta}_n = \text{argmax}_\theta \; \ell(\theta; x_{1:n})$$

$$X_1 \longrightarrow X_n \overset{iid}{\sim} P_{\theta_*}$$

$$\xrightarrow{\text{assumptions}} \quad \widehat{\theta}_n \overset{d}{\approx} N(\theta_*, (n \mathcal{I}_*)^{-1})$$

$$\mathcal{I}_* = \mathbb{E}_{\theta_*}[\nabla_\theta^2 - \log P_\theta(X)]$$

$$= \text{Cov}_{\theta_*}[\nabla_\theta^2 \log P_\theta(X)]_{\theta = \theta_*}$$

$$\approx \text{Cov}_{\widehat{\theta}_n}[\nabla_\theta^2 \log P_\theta(X)|_{\theta = \widehat{\theta}_n}]$$

$\text{don't do}$
$\pi$ with $n$?

$$\Rightarrow \text{Var}(\widehat{\theta}_{MLE}) = 1/n \cdot \mathcal{I}_* (1 + o(1))$$

Can we find $\widehat{\theta}$ s.t. $\text{Var}(\widehat{\theta}_{MLE}) \ll 1/n \mathcal{I}_*$?

Cramér-Rao : Let $\widehat{\theta}$ be unbiased for $\theta$, ie

$$\forall \theta \subset \Theta, \quad \int P_\theta(x) \widehat{\theta}(x) \, dx = \theta.$$

Then $\text{Cov}_\theta(\widehat{\theta}(x)) \geq \mathcal{I}(\theta)^{-1}$.

so, $\forall v, \; \text{Cov}_\theta(\langle \widehat{\theta}(x), v \rangle) \geq v^T \mathcal{I}(\theta)^{-1} v \in \mathbb{R}$

(1)

Intuition: If $\mathcal{E}(P, f) \gtrsim \gamma |f|^2$,

then $P \underset{\sim}{\sim} \gamma$ independent samples from $\pi$.

decorrelation time $\approx \gamma^{-1}$ of equilibrium

Comparison: If $\mathcal{E}(P_1, f) \geq K \mathcal{E}(P_2, f)$

then $P_1 \underset{\approx}{\sim} K$ steps of $P_2$

We will follow in this direction.

$P_1$ = implementable algorithm
$P_2$ = ideal slice sampling

(actually: more generally, $|f|^2 \leq S \cdot \mathcal{E}(P, f) + \beta(s) |f|^2_\infty$

even if $\beta(s) > 0$, (non-exp) convergence, e.g.

$$\| P^n f \|^2_2 \leq \gamma(n) \cdot \| f \|^2_\infty \qquad \underline{\text{Examples}}$$

so, we might instead prove

$$\mathcal{E}(P_2, f) \leq S \cdot \mathcal{E}(P_1, f) + \beta(s) \cdot \| f \|^2_\infty$$

$\approx P_1$ not much worse than $P_2$ (depending on $\beta$)

Let $H_t$ be $\mu_t$-rev, pos. . $m(t) = \nu(G(t))$

$$\forall t, \forall s \qquad \| f \|^2_{0, \mu_t} \leq S \cdot \mathcal{E}(H_t, f) + \beta(t, s) \| f \|^2_{osc}$$

Define $\beta(s) = \int_0^{\| w \|_\infty} \beta(s, t) \, m(t) \, dt$

$$\implies \quad \boxed{\phantom{XXXX}} \leq S \cdot \mathcal{E}(H, f) + \beta(s) \| f \|^2_{osc}$$

$\mathcal{E}(U, f)$

$$\implies \forall v \in \mathbb{R}^p, \quad v^T(\mathcal{I}(\theta) - C(\theta)^{-1})v \geq 0$$
$$\mathcal{I}(\theta) \succeq C(\theta)^{-1}$$
$$\underbrace{C(\theta)}_{\substack{\text{variance} \\ \text{of unbiased} \\ \text{estimator}}} \succeq \underbrace{\mathcal{I}(\theta)^{-1}}_{\substack{\text{"Cramér-Rao} \\ \text{lower bound"}}}$$

Current reaction: is this lower bound achievable?

<u>roughly</u>: asymptotically (large-n),
for reasonable models, yes,
and the MLE does so.

The MLE is <u>almost unbiased</u> (as $n \to \infty$),
<u>almost minimal-variance</u> (w.r.t. CRLB)
$\implies$ "optimal".

<u>Caveats</u>: • At {finite n ✓ large p / ...}, more to life
than ~~unbiasedness~~ unbiasedness, variance.
• All relies on correct model specification.
If $P_* \notin \{P_\theta : \theta \in \omega\}$, no luck
$\to$ doesn't say much about robustness.
• Sharp when $\eta(x)$ is parallel to $\nabla_\theta \log P_\theta(x)$.
$\to$ exponential families (can develop)
• needs "regular" model (no unif. ...)

$\Rightarrow$ ① study ~~$[\ldots]$~~ $u^n \longrightarrow \pi$

② study $H_t^n \longrightarrow \nu_t$

$\Longrightarrow$ study $H \longrightarrow \pi$.

Is the ideal slice sampler good?
Is the approximation good?

Metropolis $(\pi, \nu, \alpha)$., $\nu \otimes \alpha$ symmetric

$\begin{cases} \text{RWM}: (\pi, \text{Leb}, \text{CRW}) \\ \text{IMH} \quad (\pi, q, q) \\ \text{pCN} \quad (\pi, \gamma_{m,c}, \text{OU}) \end{cases}$

Metropolis $\subseteq$ HybridSlice
Metropolis $\lesssim$ HybridSlice

$\quad$ RWM $\lesssim$ Slice $(\pi, \text{Leb})$
$\quad$ IMH $\lesssim$ Slice $(\pi, q)$
$\quad$ pCN $\lesssim$ Slice $(\pi, \gamma_{m,c})$

Stepping-Out + Shrinkage $\qquad \leftarrow h \rightarrow$ ?

$\mathcal{E}(H_t, f) \geq \lambda(t) \|f\|_{\nu_t}^2$

$\lambda(t) = \dfrac{h - \delta(t)}{h} \times \dfrac{m(t)}{m(t) + \delta(t)}$

$\geq \dfrac{h - \Delta}{h} \times \dfrac{m(t)}{m(t) + \Delta}$

$\Rightarrow \mathcal{E}(H, f) \geq \dfrac{h - \Delta}{h} \times \dfrac{m_-}{m_- + \Delta} \mathcal{E}(u, f)$.

$h = 2\Delta \Rightarrow \geq \dfrac{1}{2} \dfrac{m_-}{m_- + \Delta}$

So, $\mathbb{E}_\theta [\nabla_\theta \log P_\theta(x)] = 0$

$I(\theta) = \mathbb{E}_\theta [\nabla_\theta \log P_\theta(x) \cdot \nabla_\theta \log P_\theta(x)^\top]$

$\qquad = \text{Cov}_\theta [\nabla_\theta \log P_\theta(x)]$

$\qquad = \mathbb{E}_\theta [\nabla_\theta^2 (-\log P_\theta(x))]$  (do examples!)

One move:

③ Since $\int P_\theta(x) \eta(x) \, dx = \theta$  (by assumption)

$\frac{d}{d\theta} \Rightarrow \int P_\theta(x) \eta(x) \nabla_\theta \log P_\theta(x) dx = I_p$

$\mathbb{E}_\theta [\underbrace{\eta(X)}_{\text{mean } \theta} \underbrace{\nabla_\theta \log P_\theta(X)^\top}_{\text{mean } 0}] = I_p$

$\text{Cov}_\theta [\eta(X), \nabla_\theta \log P_\theta(X)] = I_p$

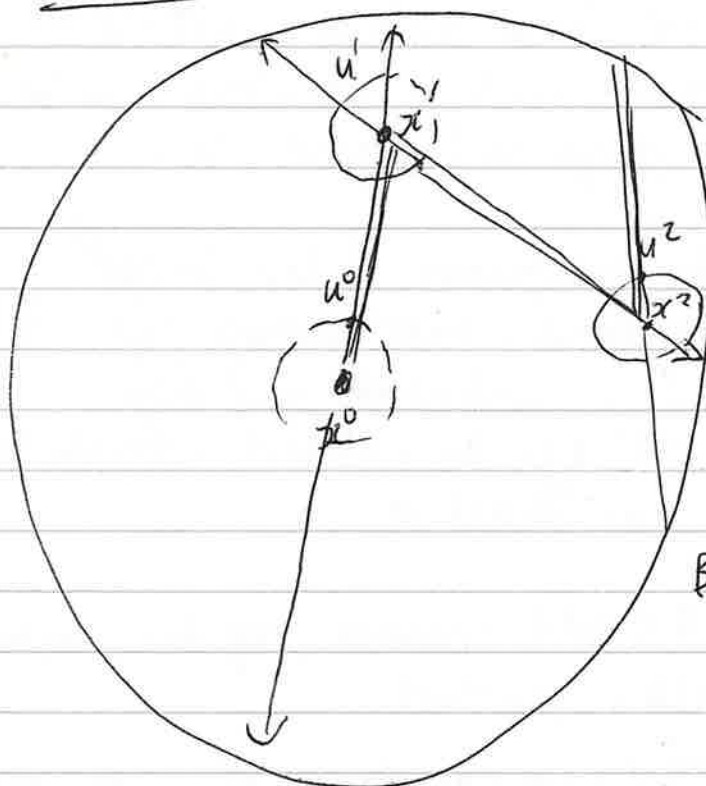Now, let $Y = \begin{pmatrix} \eta(X) \\ \nabla_\theta \log P_\theta(X) \end{pmatrix}$.

Compute $\text{Cov}_\theta(Y) = \begin{pmatrix} C(\theta) & I_p \\ I_p & \mathcal{I}(\theta) \end{pmatrix} \succeq 0.$

$\begin{pmatrix} u \\ v \end{pmatrix}^\top \text{Cov}_\theta(Y) \begin{pmatrix} u \\ v \end{pmatrix} = u^\top C(\theta) u + 2u^\top v + v^\top \mathcal{I}(\theta) v \geq 0$

min wrt $u \Rightarrow u^* = -C(\theta)^{-1} v$

$\Rightarrow \left(\underset{u}{-C(\theta)^{-1} v}\right)^\top \text{Cov}_\theta(Y) \left(\underset{v}{-C(\theta)^{-1} v}\right) = v^\top \left\{ \mathcal{I}(\theta) - C(\theta)^{-1} \right\} v$

2

Hit-and-Run     Unif($K$) , also $\sigma^{th}$-order.



$B(x,r) \supseteq K \supseteq B(x,r)$

LV: Hit-and-Run has $\gamma \gtrsim 2^{-33} \cdot \frac{1}{d^2} \left( \frac{r_K}{R_K} \right)^2$

dimension / conditioning

$\Longrightarrow$ Slice Sampling w/ quasi-concave $\pi$.

$\underline{\text{П38}}$ $m \le V'' \le L \Rightarrow \mathcal{E}(U, f) \le 2^{33} \cdot d^2 \cdot \left( \frac{L}{m} \right) \mathcal{E}(H, f)$.

$\underline{\text{П39}}$ $\pi_{d,m} \Rightarrow \| f \|^2_{d,m} \le 2^{33} \cdot d^2 \cdot \frac{(d+1)(d+m-1)}{m-1} \mathcal{E}(H, f)$

$\underline{\text{П40}}$ : $U(x) \in \| x \|^{p_1, p_2} \not\to 0, \| x \|^{q_1, q_2} \not\to \infty$

$\Rightarrow k_q(t) \lesssim \left( \log \frac{1}{t} \right)^{\frac{1}{q_1} - \frac{1}{q_2}} \cdot 0^+$

$\left( \log \left( \frac{1}{t} \right) \right)^{\frac{1}{p_1} - \frac{1}{p_2}} 1^-$

$\Longrightarrow p_1 = p_2 , \quad \beta(s) \lesssim e^{-\Omega(s^{q_1 q_2 (d s_2 - s_1)})}$

$\quad p_1 \ne p_2 \quad \beta(s) \lesssim s^{-(1 + \frac{d}{q_1}) \frac{p_1 p_2}{p_1 - p_2}}$

$\textcircled{\scriptsize 10}$

# Parametric Estimation and the CRLB

OLS : model $\to$ linear, unbiased estimators
Gauss-Markov $\to$ OLS is minimum variance
not whole story, but a good start

general parametric estimation : too broad?
· finite-dimensional, fixed sample size $\{P_\theta(x)\}$
· "linear" = ? , unbiased : still busy
  $\underline{Q}$ : how "good" can an unbiased estimator be?
  $\underline{A}$ : CRLB: $\text{Var}(\hat\theta(x)) \geq \cdots$

setting $\{P_\theta(x) : \theta \in \Theta\}$ smooth, regular. may have $x = (x_1 : x_n)$
let $\hat\theta(X) = \eta(X)$ , so that

$$\forall \theta \in \Theta, \quad \int P_\theta(x) \eta(x) \, dx = \theta \qquad \text{as vectors}$$

Measure uncertainty via

$$\text{Cov}_\theta(\hat\theta(x)) = \int P_\theta(x) (\eta(x) - \theta)(\eta(x) - \theta)^T = C(\theta).$$

What can we say about $C(\theta)$?

Some basic results $\qquad \left( \frac{d}{d\theta} f(\theta) = f(\theta) \cdot \frac{d}{d\theta} \log f(\theta) \right)$

① since $\int P_\theta(x) \, dx = 1$,
$\frac{d}{d\theta} \Rightarrow \quad \int P_\theta(x) \nabla_\theta \log P_\theta(x) \, dx = 0_p$
$\qquad \qquad \mathbb{E}_\theta[\nabla_\theta \log P_\theta(x)] = 0_p$

② since $\int P_\theta(x) \, dx = 1$
$(\frac{d}{d\theta})^2 \Rightarrow \int P_\theta(x) \{ \nabla_\theta \log P_\theta(x) \nabla_\theta \log P_\theta(x)^T + \nabla_\theta^2 \log P_\theta(x) \} dx$
$\qquad \qquad \qquad = 0_{p \times p}$

$\mathcal{I}(\theta) := \int P_\theta(x) (\nabla_\theta \log P_\theta(x))(\nabla_\theta \log P_\theta(x))^T \, dx$
$\qquad = \int P_\theta(x) (\nabla_\theta^2 (- \log P_\theta(x)) \, dx$

1

# Takeaways

- SS : good in theory and practice

- HSS : used in practice, gap in theory

- Comparison : how well does $H$ approximate $U$?

- Applications : $H \gtrsim U$

    If $H_f \succeq c$, then $H \succeq cU$

    If $H_f$ good, then $H$ almost as good as $U$.

- We focus on $H_f$ = Hit-and-Run, and similar.
  Easy to combine bounds.

- General comparison framework valid beyond SS

# Maximum Likelihood Inference

Linear Model: interpretable, fast/direct/transparent, flexible, exact inference/pivots/conjugacy

But, not always appropriate for problems/data/..

## Extension 1   Linear Mixed Models.

$m$ subjects, $n$ input-output pairs per subject

for $j = 1, -, m$,    $y_{ij} = \beta_0 + \beta_1 x_{ij} + \epsilon_{ij}$

or: subject-wise intercept $\Rightarrow$ $y_{ij} = \beta_0 + \beta_1 x_{ij} + \sum I[j=k] \gamma_h + \epsilon_j$

Can't estimate $\gamma_j$ well if $n$ is small (not identified)

shrinkage/sharing of information: $(\gamma_1, ..., \gamma_m) \sim N(0, \Sigma(\phi))$

(~ identifiability, regularisation)    e.g. $\phi I$.

## Extension 2   Non-Gaussian Input-Output Regression Model

count data: $y_j = Poisson(\lambda_j)$

$\log \lambda_j = \beta_0 + \beta_1 x_j$    ($\Rightarrow$ (latent) linear ~)

estimation by LS inappropriate $\rightarrow$ $\log$ nonlinear, $Var(y|x) $ nonconstant

· instead of forcing problems into LS format, identify salient features of LS problem which are structurally important, beyond "full tractability"

· What principles do we have? how do we come up with OLS?