

Background: Two key problems: Estimation of Integrals, Sampling

$$I = \int_{\mathbb{X}} f(x) dx$$

Reason: Want to learn , often rely on samples to approximate

Basic MC gets $\hat{I}_n = \frac{1}{n} \sum_{i=1}^n \varphi(\theta_i)$, when $I = \int_{\Theta} \varphi(\theta) \pi(\theta) d\theta$, pi pdf and psi 'test' function

Proposition 2.1 (LLN): If $\mathbb{E}(|\varphi(X)|) < \infty$ then \hat{I}_n is a strongly consistent estimator (of I)

Proposition 2.2 (CLT): If $\sigma^2 = \mathbb{V}(\varphi(X)) = \int [\varphi(x) - I]^2 \pi(x) dx < \infty$ then

$$\mathbb{E}((\hat{I}_n - I)^2) = \mathbb{V}(\hat{I}_n) = \frac{\sigma^2}{n} \quad \text{and} \quad \frac{\sqrt{n}}{\sigma} (\hat{I}_n - I) \xrightarrow{D} \mathcal{N}(0, 1)$$

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (\varphi(X_i) - \hat{I}_n)^2$$

Proposition 2.3: If conditions for 2.2 is satisfied then is unbiased for sigma^2

Proof: let $Y_i = \varphi(X_i)$ write $E[S_n^2]$ in terms of this, get $E[Y_{bar}^2]$ and sub in

$$\mathbb{P}\left(\left|\hat{I}_n - I\right| > c \frac{\sigma}{\sqrt{n}}\right) \leq \frac{\mathbb{V}(\hat{I}_n)}{c^2 \sigma^2 / n} = \frac{1}{c^2}.$$

Error Estimates: Chebyshev's Inequality

CLT $\mathbb{P}\left(\left|\hat{I}_n - I\right| > c \frac{\sigma}{\sqrt{n}}\right) \approx 2(1 - \Phi(c)) = O\left(\frac{e^{-c^2/2}}{c}\right).$, CI $\left(\hat{I}_n \pm c_\alpha \frac{\sigma}{\sqrt{n}}\right) \approx \left(\hat{I}_n \pm c_\alpha \frac{S_n}{\sqrt{n}}\right)$

Toy Example: draw points randomly in a 2x2 square, the proportion of points in the circle circumscribed by the square is $\pi/4$, can use this to get estimate of pi

Drawing Random Numbers: How do we do it? Needs to be *pseudo-random*, can't do fully

One approach: $x_{n+1} = (ax_n + c) \bmod M$, maximum period M (then cycles). Looks quite uniform and random

Assuming we can sample from uniform[0,1]: Galton's machine to draw normal samples (the physical machine



, or **Inversion Method**; define **generalised inverse**

$$F^-(u) = \inf \{x \in \mathbb{R}; F(x) \geq u\}$$

, F is cdf, F- is inverse (can do as F is right cont)

Proposition 4.1: Let F be a cdf and $U \sim \mathcal{U}_{[0,1]}$. Then $X = F^-(U)$ has cdf F.

Fact: $F^-(u) \leq x \Leftrightarrow u \leq F(x)$. Thus for $U \sim \mathcal{U}_{[0,1]}$, we have $\mathbb{P}(F^-(U) \leq x) = \mathbb{P}(U \leq F(x)) = F(x)$.

Lectures 3&4 are always assessed

Transformation method:

We can simulate $Y \sim q$, want to simulate $X \sim f(q)$, pdf = π

$$\pi(x) = q \circ \varphi^{-1}(x) \left| \det(D\varphi^{-1})(x) \right|, \quad \varphi \text{ is a bijection,}$$

Change of variables formula gets us

Note: Sometimes will be easier to do the $^{-1}$ outside the determinant

E.g. getting Gamma(a,b) from exp(1)'s. Proof by checking MGFs, *checking MGFs is a common strategy that works*

Box-Muller Algorithm - For simulating Gaussian Distribution

$$U_1, U_2 \text{ uniform, } R = \sqrt{-2 \log(U_1)}, \quad \vartheta = 2\pi U_2, \quad R^2 \sim \text{Exp}(1/2), \quad \vartheta \sim \mathcal{U}_{[0, 2\pi]},$$

$X = R \cos(\vartheta), Y = R \sin(\vartheta)$ a bijection, Get $f(X, Y)$ in terms of $f(R^2, \theta)$ using change of variables, find (X, Y) are independent standard normal

Also a method for **multivariate normal** in slides 3, similar only difference is this time we get the inverse on the inside

$$\pi(x) = \int \bar{\pi}_{X,Y}(x,y) dy$$

Sampling via composition: Assume π marginal, i.e.

$$\bar{\pi}_{X,Y}(x,y) = \bar{\pi}_Y(y) \bar{\pi}_{X|Y}(x|y).$$

Sometimes it's easy to compute π but hard/impossible to get π , so can sample $Y \sim \bar{\pi}_Y$ then $X|Y \sim \bar{\pi}_{X|Y}(\cdot|Y)$ so $(X, Y) \sim \bar{\pi}_{X,Y}$ and hence $X \sim \pi$.

Rejection Sampling 1. Draw $X \sim q$, draw $U \sim \mathcal{U}_{[0,1]}$. $\pi(x) \leq \tilde{M} q(x)$

2. Accept $X = x$ as a sample from π if $U \leq \frac{\pi(x)}{\tilde{M} q(x)}$, constraint need such a π to exist

Proposition: The distribution of the samples accepted by rejection sampling is π . Proof in slides, uses Bayes on $\mathbb{P}(X \in A | X \text{ accepted})$, writes being in A and constrain on U as indicator functions to get that probability under $q(x)$

$$\begin{aligned} \mathbb{P}(X \in A | X \text{ accepted}) &= \frac{\mathbb{P}(X \in A, X \text{ accepted})}{\mathbb{P}(X \text{ accepted})} \quad \mathbb{P}(X \in A, X \text{ accepted}) \\ &= \int_{\mathbb{X}} \int_0^1 \mathbb{I}_A(x) \mathbb{I}\left(u \leq \frac{\pi(x)}{\tilde{M} q(x)}\right) q(x) du dx = \int_{\mathbb{X}} \mathbb{I}_A(x) \frac{\pi(x)}{\tilde{M} q(x)} q(x) dx = \int_{\mathbb{X}} \mathbb{I}_A(x) \frac{\pi(x)}{\tilde{M}} dx = \frac{\pi(A)}{\tilde{M}}. \end{aligned}$$

Often we only know π and q up to some unknown normalising constant $\pi = \tilde{\pi}/Z_\pi$ and $q = \tilde{q}/Z_q$

If you can upper bound $\tilde{\pi}(x)/\tilde{q}(x) \leq \tilde{M}$, then using $\tilde{\pi}$, \tilde{q} and \tilde{M} in the algorithm is correct.

$$\frac{\tilde{\pi}(x)}{\tilde{q}(x)} \leq \tilde{M} \Leftrightarrow \frac{\pi(x)}{q(x)} \leq \tilde{M} \frac{Z_q}{Z_\pi} = M.$$

Lemma 1: Let T denote the number of pairs (X, U) that have to be generated until X is accepted for the first time. T is geometrically distributed with parameter $1/M$ and in particular $E(T) = M$. Gives a way to estimating ratio of

$$\mathbb{E}(T) = M = \frac{Z_q \tilde{M}}{Z_\pi}$$

normalising constants as Example: Uniform from bounded subset B of R^p , set

$$\tilde{M} = 1 \quad (\text{If the bounded subset filled the uniform}), \quad \tilde{\pi}(x) / (\tilde{M} \tilde{q}(x)) = \mathbb{I}_B(x), \quad \text{prob(accept)} = \frac{Z_\pi}{Z_q}.$$

In summary: Requires evaluation of π pointwise, an upper bound of π/q (or same on tildas). Upper Bound not always feasible

Importance Sampling: Want to compute $I = \mathbb{E}_\pi(\varphi(X)) = \int_{\mathbb{X}} \varphi(x) \pi(x) dx$, don't know π but have a proposal distribution q

Require the support of π to be in q $\pi(x) > 0 \Rightarrow q(x) > 0$; q called the **proposal/importance distribution**

Identity $I = \mathbb{E}_\pi(\varphi(X)) = \mathbb{E}_q(\varphi(X)w(X))$. w is the **importance weight function** $\pi/q \Rightarrow$ for X iid q ,

$$\hat{I}_n^{IS} = \frac{1}{n} \sum_{i=1}^n \varphi(X_i) w(X_i).$$

properties of importance sampling: Unbiased, strongly consistent, CLT

$$\mathbb{V}_q(\hat{I}_n^{IS}) = \sigma_{IS}^2 / n \quad \sigma_{IS}^2 := \mathbb{V}_q(\varphi(X)w(X)) \quad \text{If } \sigma_{IS}^2 < \infty \quad \lim_{n \rightarrow \infty} \sqrt{n} (\hat{I}_n^{IS} - I) \xrightarrow{D} \mathcal{N}(0, \sigma_{IS}^2)$$

Same proofs as for rejection sampling

Question: Is there a best proposal that minimises the variance sigma (IS)? **Proposition:** optimisal proposal

We have indeed

$$\sigma_{IS}^2 = \mathbb{V}_q(\varphi(X)w(X)) = \mathbb{E}_q(\varphi^2(X)w^2(X)) - I^2.$$

We also have by Jensen's inequality for any q

$$\mathbb{E}_q(\varphi^2(X)w^2(X)) \geq \left(\int_{\mathbb{X}} |\varphi(x)| \pi(x) dx \right)^2.$$

For $q = q_{opt}$, we have

$$\begin{aligned} \mathbb{E}_{q_{opt}}(\varphi^2(X)w^2(X)) &= \int_{\mathbb{X}} \frac{\varphi^2(x)\pi^2(x)}{|\varphi(x)|\pi(x)} dx \times \int_{\mathbb{X}} |\varphi(x)| \pi(x) dx \\ &= \left(\int_{\mathbb{X}} |\varphi(x)| \pi(x) dx \right)^2. \end{aligned}$$

$$q_{opt}(x) = \frac{|\varphi(x)| \pi(x)}{\int_{\mathbb{X}} |\varphi(x)| \pi(x) dx}$$

minimising $V_q(I^{nIS})$ is
(modulus to make it a true density)]

Proof Variance equation gets an equation with a part I we can't minimise on
Jensens gives us a lower bound, subbing in q_{opt}
gives us equality to the lower bound

Problem: $q_{opt}(x)$ can never be used in practice. Requires us to know I Can use as a guideline to choose q though

Not always best to sample from the actual distribution π , e.g. another distribution that is most concentrated at the most likely points of π can be better (toy example in slide 4, t distribution)

Normalised Importance Sampling: Assume $\pi(x) = \tilde{\pi}(x)/Z_\pi$ and $q(x) = \tilde{q}(x)/Z_q$,

$$\pi(x) > 0 \Rightarrow q(x) > 0 \text{ and and define } \tilde{w}(x) = \frac{\tilde{\pi}(x)}{\tilde{q}(x)} \quad \text{Alternative Identity is}$$

$$I = \mathbb{E}_\pi(\varphi(X)) = \frac{\int_{\mathbb{X}} \varphi(x) \tilde{w}(x) q(x) dx}{\int_{\mathbb{X}} \tilde{w}(x) q(x) dx}.$$

Main complication is this is the ratio of expectations ($\pi \sim q$)

Propositions 2.1 (SLLN for NIS) Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} q$ and assume that $\mathbb{E}_q(|\varphi(X)|w(X)) < \infty$.

then is strongly consistent. **Proof:** divide top and bottom by n , both converge a.s by SLLN. BUT

for finite n \hat{I}_n^{NIS} is **biased**

Proposition 2.2 If $\mathbb{V}_q(\varphi(X)w(X)) < \infty$ and $\mathbb{V}_q(w(X)) < \infty$ then

$$\sqrt{n}(\hat{I}_n^{NIS} - I) \Rightarrow \mathcal{N}(0, \sigma_{NIS}^2), \quad \sigma_{NIS}^2 := \mathbb{V}_q([w(X)(\varphi(X) - I)]) = \int \frac{\pi(x)^2 (\varphi(x) - I)^2}{q(x)} dx.$$

Proof

First notice that with X_1, \dots, X_n i.i.d. $\sim q$

$$\sqrt{n}(\hat{I}_n^{\text{NIS}} - I) = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{w}(X_i)[\varphi(X_i) - I]}{\frac{1}{n} \sum_{i=1}^n \tilde{w}(X_i)}$$

where since $\tilde{w}(x) = \tilde{\pi}/\tilde{q}$

$$\mathbb{E}_q[\tilde{w}(X_n)(\varphi(X_i) - I)] = 0.$$

Since $\mathbb{V}_q(\varphi(X)w(X)) < \infty$ by standard CLT

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{w}(X_i)[\varphi(X_i) - I] \Rightarrow \mathcal{N}\left(0, \mathbb{V}_q(\tilde{w}(X_1)[\varphi(X_1) - I])\right).$$

The strong law of large numbers applied to the denominator

$$\frac{1}{n} \sum_{i=1}^n \tilde{w}(X_i) \rightarrow \mathbb{E}_q[\tilde{w}(X_1)] = Z_\pi/Z_q, \quad \text{a.s.}$$

By Slutsky's theorem, combining the two

$$\begin{aligned} \sqrt{n}(\hat{I}_n^{\text{NIS}} - I) &\Rightarrow \mathcal{N}\left(0, \mathbb{V}_q(\tilde{w}(X_1)[\varphi(X_1) - I]) \frac{Z_q^2}{Z_\pi^2}\right) \\ &\sim \mathcal{N}\left(0, \sigma_{\text{NIS}}^2\right). \end{aligned}$$

Variance of importance Sampling estimators

- Standard Importance Sampling: $X_1, \dots, X_n \stackrel{iid}{\sim} q$,

$$\hat{I}_n^{\text{S}} = \frac{1}{n} \sum_{i=1}^n \varphi(X_i)w(X_i).$$

- Asymptotic Variance:

$$\begin{aligned} \mathbb{V}_{as}(\hat{I}_n^{\text{S}}) &= \mathbb{E}_q[(\varphi(X)w(X) - \mathbb{E}_q(\varphi(X)w(X)))^2] \\ &\approx \frac{1}{n} \sum_{i=1}^n (\varphi(X_i)w(X_i) - \hat{I}_n^{\text{S}})^2. \end{aligned}$$

- Thus the asymptotic variance can be estimated consistently with

$$\frac{1}{n} \sum_{i=1}^n (\varphi(X_i)w(X_i) - \hat{I}_n^{\text{S}})^2.$$

- Normalised Importance Sampling: $X_1, \dots, X_n \stackrel{iid}{\sim} q$,

$$\hat{I}_n^{\text{NIS}} = \frac{\sum_{i=1}^n \varphi(X_i)\tilde{w}(X_i)}{\sum_{i=1}^n \tilde{w}(X_i)}.$$

- Asymptotic Variance:

$$\mathbb{V}_{as}(\hat{I}_n^{\text{NIS}}) = \frac{\mathbb{E}_q[(\tilde{w}(x)(\varphi(X) - I))^2]}{\mathbb{E}_q[\tilde{w}(X)]^2}.$$

- Thus the asymptotic variance can be estimated consistently with

$$\frac{\frac{1}{n} \sum_{i=1}^n \tilde{w}(X_i)^2 (\varphi(X_i) - \hat{I}_n^{\text{NIS}})^2}{\left(\frac{1}{n} \sum_{i=1}^n \tilde{w}(X_i)\right)^2}.$$

MC vs IS vs NIS: NIS is sometimes better than IS, despite seeming to have extra variance from the ratio of distributions and biased (while IS is unbiased)

- **Toy example:** Let $\mathbb{X} = \mathbb{R}^d$ and

$$\pi(x) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\sum_{i=1}^d x_i^2}{2}\right)$$

and

$$q(x) = \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{\sum_{i=1}^d x_i^2}{2\sigma^2}\right).$$

Important toy example

$$\mathbb{P}(X \text{ accepted}) = \frac{1}{\sigma^d} \rightarrow 0 \text{ as } d \rightarrow \infty,$$

very inefficient as d gets big and $\sigma > 1$

Writes out LHS., take expectation, and then use CLT

Use SLLN, then **Slutskys Theorem** to combine

(Alternatively use Delta method instead of Slutskys)

SLLN, CLT, and Slutskys are the three main tools we will use'

$$\mathbb{V}_q [w(X)] = \left(\frac{\sigma^4}{2\sigma^2 - 1} \right)^{d/2} - 1$$

Importance sampling, exponentially increasing variance

Weird as we wanted the rate to be independent of d, and here this is true but the ‘constant’ (in n) explodes exponentially in d

Markov Chain Monte Carlo:

Pseudo Marginal is *not examinable*

Markov Chain Theory is *not examinable*, application is (won’t get question like what is irreducible)

Hamiltonian Monte Carlo is *not examinable kinda*, too hard to make a question on

1 question from MCMC, Gibbs Sampling, Metropolis Hasting, Parallel Tempering Capacity and Diagnostics are examinable (PTC and Diagnostics are part of Metropolis Hasting kinda)

MC methods give convergence to pi of rate 1/Root(n), independent of dimension d of x. But error still depends on d, through constant in front of rate, and it usually explodes in d. *If we could get something dependence of $\pi(x)/\pi(y)$ this constant would disappear*

$$\frac{1}{n} \sum_{t=1}^n \phi(X_t) \rightarrow \int \phi(x) \pi(x) dx$$

Want a distribution X such that works for X correctly distributed by Ergodic theorem, CLT for rate but we have a Markov Chain, so don’t have independence, so need a different thing to get the rate.

- ▶ The state space \mathbb{X} is now continuous, e.g. \mathbb{R}^d .
- ▶ $(X_t)_{t \geq 1}$ is a Markov chain if for any (measurable) set A ,

$$\begin{aligned} \mathbb{P}(X_t \in A | X_1 = x_1, X_2 = x_2, \dots, X_{t-1} = x_{t-1}) \\ = \mathbb{P}(X_t \in A | X_{t-1} = x_{t-1}). \end{aligned}$$

- ▶ We have

$$\mathbb{P}(X_t \in A | X_{t-1} = x) = \int_A K(x, y) dy = K(x, A),$$

that is conditional on $X_{t-1} = x$, X_t is a random variable which admits a probability density function $K(x, \cdot)$.

- ▶ $K : \mathbb{X}^2 \rightarrow \mathbb{R}$ is the **kernel** of the Markov chain.

► Consider the autoregressive (AR) model

$$X_t = \rho X_{t-1} + V_t$$

where $V_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \tau^2)$. This defines a Markov chain such that

$$K(x, y) = \frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{1}{2\tau^2}(y - \rho x)^2\right).$$

Autoregressive model

Continuous markov chain definitions

reducibility and aperiodicity

Definition

Given a probability measure μ over \mathbb{X} , a Markov chain is μ -irreducible if

$$\forall x \in \mathbb{X} \quad \forall A : \mu(A) > 0 \quad \exists t \in \mathbb{N} \quad K^t(x, A) > 0.$$

A μ -irreducible Markov chain of transition kernel K is **periodic** if there exists some partition of the state space $\mathbb{X}_1, \dots, \mathbb{X}_d$ for $d \geq 2$, such that

$$\forall i, j, s : \mathbb{P}(X_{t+s} \in \mathbb{X}_j | X_t \in \mathbb{X}_i) = \begin{cases} 1 & j = i + s \bmod d \\ 0 & \text{otherwise.} \end{cases}.$$

Otherwise the chain is **aperiodic**.

Recurrence and Harris Recurrence

For any measurable set A of \mathbb{X} , let

$$\eta_A = \sum_{k=1}^{\infty} \mathbb{1}_A(X_k),$$

be the number of visits to the set A .

Definition

A μ -irreducible Markov chain is **recurrent** if for any measurable set $A \subset \mathbb{X} : \mu(A) > 0$, then

$$\forall x \in A \quad \mathbb{E}_x(\eta_A) = \infty.$$

A μ -irreducible Markov chain is **Harris recurrent** if for any measurable set $A \subset \mathbb{X} : \mu(A) > 0$, then

$$\forall x \in \mathbb{X} \quad \mathbb{P}_x(\eta_A = \infty) = 1.$$

Note: Harris recurrence is stronger than recurrence.

Invariant Distribution and Reversibility

Definition

A distribution of density π is **invariant** or **stationary** for a Markov kernel K , if

$$\int_{\mathbb{X}} \pi(x) K(x, y) dx = \pi(y).$$

A Markov kernel K is π -reversible if

$$\begin{aligned} \forall f \quad & \iint f(x, y) \pi(x) K(x, y) dx dy \\ &= \iint f(x, y) \pi(y) K(y, x) dx dy \end{aligned}$$

where f is a bounded measurable function.

Gibbs Sampling

$$\pi(x) = \pi(x_1, x_2, \dots, x_d), \quad x \in \mathbb{R}^d.$$

Interested in sampling

Notation: $x_{-i} := (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_d)$.

Systematic scan Gibbs sampler. Let $(X_1^{(1)}, \dots, X_d^{(1)})$ be the initial state then iterate for $t = 2, 3, \dots$

1. Sample $X_1^{(t)} \sim \pi_{X_1|X_{-1}}(\cdot | X_2^{(t-1)}, \dots, X_d^{(t-1)})$.
⋮
- j. Sample $X_j^{(t)} \sim \pi_{X_j|X_{-j}}(\cdot | X_1^{(t)}, \dots, X_{j-1}^{(t)}, X_{j+1}^{(t-1)}, \dots, X_d^{(t-1)})$.
⋮
- d. Sample $X_d^{(t)} \sim \pi_{X_d|X_{-d}}(\cdot | X_1^{(t)}, \dots, X_{d-1}^{(t)})$.

This is a markov chain: New state only depends on the current state. If instead of using the updated $X_j(t)$ s in updating future ones and just used $X_j(t-1)$ we'd still have a markov chain, but a different one
A few questions about this algorithm:

1. Is the joint distribution π uniquely specified by the conditional distributions $\pi_{X_i|X_{-i}}$?
2. Does the Gibbs sampler provide a Markov chain with the correct stationary distribution π ?
3. If yes, does the Markov chain converge towards this invariant distribution?

It will turn out to be the case under some mild conditions

Hammersley–Clifford Theorem

Consider a distribution with continuous density $\pi(x_1, x_2, \dots, x_d)$
such that

$$supp(\pi) = supp\left(\bigotimes_{i=1}^d \pi_{X_i}\right).$$

Then for any $(z_1, \dots, z_d) \in supp(\pi)$, we have

$$\pi(x_1, x_2, \dots, x_d) \propto \prod_{j=1}^d \frac{\pi_{X_j|X_{-j}}(x_j | x_{1:j-1}, z_{j+1:d})}{\pi_{X_j|X_{-j}}(z_j | x_{1:j-1}, z_{j+1:d})}.$$

Note: the condition above is known as the **positivity condition**.

Equivalently, if $\pi_{X_i}(x_i) > 0$ for $i = 1, \dots, d$, then

$$\pi(x_1, \dots, x_d) > 0.$$

Proof.

We have

$$\begin{aligned} \pi(x_{1:d-1}, x_d) &= \pi_{X_d|X_{-d}}(x_d | x_{1:d-1})\pi(x_{1:d-1}), \\ \pi(x_{1:d-1}, z_d) &= \pi_{X_d|X_{-d}}(z_d | x_{1:d-1})\pi(x_{1:d-1}). \end{aligned}$$

Therefore

$$\begin{aligned} \pi(x_{1:d}) &= \pi(x_{1:d-1}, z_d) \frac{\pi(x_{1:d-1}, x_d)}{\pi(x_{1:d-1}, z_d)} \\ &= \pi(x_{1:d-1}, z_d) \frac{\pi(x_{1:d-1}, x_d)/\pi(x_{1:d-1})}{\pi(x_{1:d-1}, z_d)/\pi(x_{1:d-1})} \\ &= \pi(x_{1:d-1}, z_d) \frac{\pi_{X_d|X_{-d}}(x_d | x_{1:d-1})}{\pi_{X_d|X_{-d}}(z_d | x_{1:d-1})}. \end{aligned}$$

Proof.

Similarly, we have

$$\begin{aligned} \pi(x_{1:d-1}, z_d) &= \pi(x_{1:d-2}, z_{d-1}, z_d) \frac{\pi(x_{1:d-1}, z_d)}{\pi(x_{1:d-2}, z_{d-1}, z_d)} \\ &= \pi(x_{1:d-2}, z_{d-1}, z_d) \frac{\pi(x_{1:d-1}, z_d)/\pi(x_{1:d-2}, z_d)}{\pi(x_{1:d-2}, z_{d-1}, z_d)/\pi(x_{1:d-2}, z_d)} \\ &= \pi(x_{1:d-2}, z_{d-1}, z_d) \frac{\pi_{X_{d-1}|X_{-(d-1)}}(x_{d-1} | x_{1:d-2}, z_d)}{\pi_{X_{d-1}|X_{-(d-1)}}(z_{d-1} | x_{1:d-2}, z_d)} \end{aligned}$$

hence

$$\begin{aligned} \pi(x_{1:d}) &= \pi(x_{1:d-2}, z_{d-1}, z_d) \frac{\pi_{X_{d-1}|X_{-(d-1)}}(x_{d-1} | x_{1:d-2}, z_d)}{\pi_{X_{d-1}|X_{-(d-1)}}(z_{d-1} | x_{1:d-2}, z_d)} \\ &\quad \times \frac{\pi_{X_d|X_{-d}}(x_d | x_{1:d-1})}{\pi_{X_d|X_{-d}}(z_d | x_{1:d-1})} \end{aligned}$$

Proof.

By $z \in supp(\pi)$ we have that $\pi_{X_i}(z_i) > 0$ for all i .

Also, we are allowed to suppose that $\pi_{X_i}(x_i) > 0$ for all i .

Thus all the conditional probabilities we introduce are positive since

$$\begin{aligned} &\pi_{X_j|X_{-j}}(x_j | x_1, \dots, x_{j-1}, z_{j+1}, \dots, z_d) \\ &= \frac{\pi(x_1, \dots, x_{j-1}, x_j, z_{j+1}, \dots, z_d)}{\pi(x_1, \dots, x_{j-1}, z_{j+1}, \dots, z_d)} > 0. \end{aligned}$$

By iterating we have the theorem. \square

Example: Non-Integrable target (in slide 6)

Invariance of the Gibbs sampler

The systematic scan Gibbs sampler kernel admits π as invariant distribution.

- The kernel of the Gibbs sampler (case $d = 2$) is

Proof for $d = 2$.

Let $x = (x_1, x_2)$ and $y = (y_1, y_2)$. Then we have

$$K(x^{(t-1)}, x^{(t)}) = \pi_{X_1|X_2}(x_1^{(t)} | x_2^{(t-1)})\pi_{X_2|X_1}(x_2^{(t)} | x_1^{(t)}) \quad \int K(x, y)\pi(x)dx = \int \pi(y_2 | y_1)\pi(y_1 | x_2)\pi(x_1, x_2)dx_1 dx_2$$

- Case $d > 2$:

$$K(x^{(t-1)}, x^{(t)}) = \prod_{j=1}^d \pi_{X_j|X_{-j}}(x_j^{(t)} | x_{1:j-1}^{(t)}, x_{j+1:d}^{(t-1)}) \quad = \pi(y_2 | y_1) \int \pi(y_1 | x_2)\pi(x_2)dx_2 \\ = \pi(y_2 | y_1)\pi(y_1) = \pi(y_1, y_2) = \pi(y).$$

Can use this for bigger d by grouping them (e.g $X_1 \sim X_1$, $X_2 \sim (X_2, X_3, \dots)$)

Proposition

Assume π satisfies the positivity condition, then the Gibbs sampler yields a π -irreducible and recurrent * Markov chain.

Proof.

Recurrence follows from irreducibility and the fact that π is invariant (see Meyn and Tweedie, Prop'n 10.1.1.).

Irreducibility. Let $\mathbb{X} \subset \mathbb{R}^d$, such that $\pi(\mathbb{X}) = 1$. Write K for the kernel and let $A \subset \mathbb{X}$ such that $\pi(A) > 0$. Then for any $x \in \mathbb{X}$

$$\begin{aligned} K(x, A) &= \int_A K(x, y)dy \\ &= \int_A \pi_{X_1|X_{-1}}(y_1 | x_2, \dots, x_d) \times \dots \\ &\quad \times \pi_{X_d|X_{-d}}(y_d | y_1, \dots, y_{d-1})dy. \end{aligned}$$

□

Irreducibility and Recurrence (cont.)

Proof.

Thus if for some $x \in \mathbb{X}$ and A with $\pi(A) > 0$ we have $K(x, A) = 0$, we must have that

$$\pi_{X_1|X_{-1}}(y_1 | x_2, \dots, x_d) \times \dots \times \pi_{X_d|X_{-d}}(y_d | y_1, \dots, y_{d-1}) = 0,$$

for almost all $y = (y_1, \dots, y_d) \in A$.

Therefore, by the Hammersley–Clifford theorem, we must also have that

$$\pi(y_1, y_2, \dots, y_d) \propto \prod_{j=1}^d \frac{\pi_{X_j|X_{-j}}(y_j | y_{1:j-1}, x_{j+1:d})}{\pi_{X_j|X_{-j}}(x_j | y_{1:j-1}, x_{j+1:d})} = 0,$$

for almost all $y = (y_1, \dots, y_d) \in A$ and thus $\pi(A) = 0$ obtaining a contradiction. □

Theorem

If the positivity condition is satisfied then for any π -integrable function $\phi : \mathbb{X} \rightarrow \mathbb{R}$:

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{i=1}^t \phi(X^{(i)}) = \int_{\mathbb{X}} \phi(x) \pi(x) dx$$

LLN for Gibbs Sampler for π -almost all starting values $X^{(1)}$.

Gibbs sampling and auxiliary variables

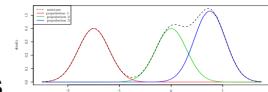
- Gibbs sampling requires sampling from $\pi_{X_j|X_{-j}}$.
- In many scenarios, we can include a set of auxiliary variables Z_1, \dots, Z_p and have an "extended" distribution of joint density $\bar{\pi}(x_1, \dots, x_d, z_1, \dots, z_p)$ such that

$$\int \bar{\pi}(x_1, \dots, x_d, z_1, \dots, z_p) dz_1 \dots dz_d = \pi(x_1, \dots, x_d). \quad \text{Mixture models,}$$

Capture-recapture,tobit,probit etc.

Sometimes adding extra 'auxillary' variables (Z_1, \dots, Z_p) actually makes the distribution easier
Can be hard to decide what auxiliary variables to get. E.g. didn't have a nice method for
sampling Bayesian Logistic Regression until 2013

Example: Mixture of Normals. $K = 2$ for simplicity, data is a mixture of gaussians with



different mean but same variance 1, want to estimate the means
Not obvious how to do it, because the conditional distribution (below) is pretty messy

- Independent data $y_1, \dots, y_n \quad Y_i|\theta \sim \sum_{k=1}^K p_k \mathcal{N}(\mu_k, \sigma_k^2)$ where $\theta = (p_1, \dots, p_K, \mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2)$

► Likelihood function

$$p(y_1, \dots, y_n|\theta) = \prod_{i=1}^n p(y_i|\theta) = \prod_{i=1}^n \left(\sum_{k=1}^K \frac{p_k}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{(y_i - \mu_k)^2}{2\sigma_k^2}\right) \right) \quad p(\theta) = p(\mu_1)p(\mu_2)$$

► For example, let's fix $K = 2$, $\sigma_k^2 = 1$ and $p_k = 1/K$ for all k . Suppose prior

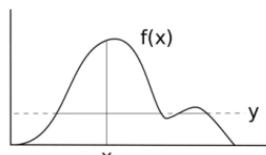
Goal: introduce a variable for each data point that tells us which gaussian the datapoint comes from, in a way that preserves original marginal density

Associate to each Y_i an auxiliary Z_i where $\mathbb{P}(Z_i = k|\theta) = p_k$ and $Y_i|Z_i = k, \theta \sim \mathcal{N}(\mu_k, \sigma_k^2)$

$$\Rightarrow p(y_i|\theta) = \sum_{k=1}^K \mathbb{P}(Z_i = k) \mathcal{N}(y_i; \mu_k, \sigma_k^2) \quad \text{extended posterior} \quad p(\theta, z_1, \dots, z_n | y_1, \dots, y_n) \propto p(\theta) \prod_{i=1}^n \mathbb{P}(z_i|\theta) p(y_i|z_i, \theta).$$

Integrating out the Z_i in the new likelihood function gives the original Likelihood

Gibbs Samples alternately $\mathbb{P}(z_{1:n}|y_{1:n}, \theta) \quad p(\theta|y_{1:n}, z_{1:n})$ (saw alternate sampling last time)
Notes has an example for set variance 1 and p_k uniform.



Slice Sampling: A form of Gibbs Sampling, y is the auxillary.

Assumes you can compute and sample the distribution along each line.

Start from a random x coordinate sample along that sample from random $y = f(x)$
proportional to $p(x)$ coordinate sample along that.

Consider the joint density over $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$ uniform
We over the region $\{(x, y) : 0 < y < f(x)\}$ $p(x, y) = (1/Z) \mathbb{1}_{\{0 < y < f(x)\}}, \quad Z = \int f(x) dx$

The marginal density for x is then $\int_0^{f(x)} (1/Z) dy = f(x)/Z = \pi(x)$

Exam question: want to show this is a gibbs sampler for a given distribution. Two steps $\pi(\cdot|y)$ and $\pi(\cdot|x)$, show they're uniform

Gibbs Sampling Summary

- ▶ Given a target $\pi(x) = \pi(x_1, x_2, \dots, x_d)$, Gibbs sampling works by sampling from $\pi_{x_j|x_{-j}}(x_j|x_{-j})$ for $j = 1, \dots, d$.
- ▶ Sampling exactly from one of these full conditionals might be a hard problem itself.
- ▶ Even if it is possible, the Gibbs sampler might converge slowly if components are highly correlated.
- ▶ If the components are not highly correlated then Gibbs sampling performs well, even when $d \rightarrow \infty$ (e.g. mixing times of some Gibbs samplers can even be dimension-free, see [Yang and Rosenthal, 2017]).
- ▶ Metropolis-Hastings algorithm (1953, 1970) is a more general algorithm that can bypass these problems.

Hard to ask a question on Gibbs because the only thing you can really do is ask us to define an auxiliary variable for something

Metropolis Hastings Algorithm: Frameworks:

1. Target distribution on $\mathbb{X} = \mathbb{R}^d$ of density $\pi(x)$
2. Proposal Distribution for any $x, x' \in \mathbb{X}$, we have $q(x'|x) \geq 0$ and $\int_{\mathbb{X}} q(x'|x) dx' = 1$.
3. Starting with $X^{(1)}$, for $t = 2, 3, \dots$
 - a. 1. Sample $X^* \sim q(\cdot|X^{(t-1)})$
 - b. Compute $\alpha(X^*|X^{(t-1)}) = \min\left(1, \frac{\pi(X^*) q(X^{(t-1)}|X^*)}{\pi(X^{(t-1)}) q(X^*|X^{(t-1)})}\right)$
 - c. Sample $U \sim U_{[0,1]}$. If $U \leq \alpha(X^*|X^{(t-1)})$, set $X^{(t)} = X^*$, otherwise set $X^{(t)} = X^{(t-1)}$.
4. Q is something we can sample from. 3c, is if we decide to accept a change to a new state, moving around the distribution
5. Is clearly a markov chain: movement only dependent on proposal and current position

$$a(x^{(t-1)}) := \int_{\mathbb{X}} \alpha(x|x^{(t-1)}) q(x|x^{(t-1)}) dx$$

Average Acceptance probability from current state

in which case $X^{(t)} = X$, otherwise $X^{(t)} = X^{(t-1)}$.

Only requires point-wise evaluations of $\pi(x)$ up to a normalising constant. if $\tilde{\pi}(x) \propto \pi(x)$

$$\frac{\pi(x^*) q(x^{(t-1)}|x^*)}{\pi(x^{(t-1)}) q(x^*|x^{(t-1)})} = \frac{\tilde{\pi}(x^*) q(x^{(t-1)}|x^*)}{\tilde{\pi}(x^{(t-1)}) q(x^*|x^{(t-1)})}.$$

Lemma: kernel of MH algorithm is $K(x, y) = \alpha(y|x)q(y|x) + (1 - a(x))\delta_x(y)$.

$$\text{Proof: } K(x, y) = \int q(x^*|x)\{\alpha(x^*|x)\delta_{x^*}(y) + (1 - \alpha(x^*|x))\delta_x(y)\}dx^*$$

$$= q(y|x)\alpha(y|x) + \left\{ \int q(x^*|x)(1 - \alpha(x^*|x))dx^* \right\} \delta_x(y) = q(y|x)\alpha(y|x) + \left\{ 1 - \int q(x^*|x)\alpha(x^*|x)dx^* \right\} \delta_x(y)$$

$$= q(y|x)\alpha(y|x) + \{1 - a(x)\} \delta_x(y).$$

Need to be careful as distribution itself is continuous, but for K can assume x is given (as for if accepted ($y=x$) or rejected ($y=x$))

Advice for above: In exam consider $y=x$ and $y=/=x$ separately

Reversibility: The Metropolis–Hastings kernel K is π -reversible and thus admits π as invariant distribution.

$$\begin{aligned}\text{Proof: For any } x, y \in \mathbb{X}, \text{ with } x \neq y \quad & \pi(x)K(x, y) = \pi(x)q(y | x)\alpha(y | x) = \pi(x)q(y | x)\left(1 \wedge \frac{\pi(y)q(x | y)}{\pi(x)q(y | x)}\right) \\ & = (\pi(x)q(y | x) \wedge \pi(y)q(x | y)) = \pi(y)q(x | y)\left(\frac{\pi(x)q(y | x)}{\pi(y)q(x | y)} \wedge 1\right) = \pi(y)K(y, x).\end{aligned}$$

If $x = y$, then obviously $\pi(x)K(x, y) = \pi(y)K(y, x)$

Consider target $\pi(x) = (\mathcal{U}_{[0,1]}(x) + \mathcal{U}_{[2,3]}(x))/2$ and proposal $q(x^* | x) = \mathcal{U}_{(x-\delta, x+\delta)}(x^*)$

The MH chain is reducible if $\delta \leq 1$: the chain stays either in $[0,1]$ or $[2,3]$

Proposition: If $q(x^* | x) > 0$ for any $x, x^* \in \text{supp}(\pi)$ then the

Metropolis-Hastings chain is **irreducible**, in fact every state can be reached in a single step (strongly irreducible).

If the MH chain is **irreducible** then it is also **Harris recurrent**

LLN for MH: If the Markov chain generated by the Metropolis–Hastings sampler is

π -irreducible, then we have for any integrable function $\phi : \mathbb{X} \rightarrow \mathbb{R}$:

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{i=1}^t \phi(X^{(i)}) = \int_{\mathbb{X}} \phi(x) \pi(x) dx \quad \text{for every starting value } X^{(1)}.$$

Independent proposal: a proposal distribution $q(x^* | x)$ not dependent on x

$$\alpha(x^* | x) = \min\left(1, \frac{\pi(x^*)q(x)}{\pi(x)q(x^*)}\right)$$

Acceptance probability then

Uniformly Ergodic: If $\pi(x)/q(x) < M$ for all x and some $M < \infty$, (hard to satisfy)

If it is satisfied, then the expected acceptance probability is at least $1/M$ (was an exam question to prove). If such M doesn't exist, not even geometrically ergodic.

Choosing a good proposal distribution:

Two sources of correlation:

Goal: design a Markov chain with small correlation $\rho(X^{(t-1)}, X^{(t)})$ between subsequent values (why?).

- ▶ between the current state $X^{(t-1)}$ and proposed value $X \sim q(\cdot | X^{(t-1)})$,
- ▶ correlation induced if $X^{(t)} = X^{(t-1)}$, if proposal is rejected.

Trade-off: there is a compromise between

- ▶ proposing large moves,
- ▶ obtaining a decent acceptance probability. Don't want to keep rejecting (so may want to be close to current place where we know we were accepted), but don't want to explore only one area of the state space

Random walk Metropolis (RWM)

In the Metropolis–Hastings, pick $q(x^* | x) = g(x^* - x)$ with g being a *symmetric* distribution, thus

$$\alpha(x^* | x) = \min\left(1, \frac{\pi(x^*)}{\pi(x)}\right)$$

So acceptance probability is

- ▶ a move to a more probable state with probability 1;
- ▶ a move to a less probable state with probability $\pi(x^*)/\pi(x) \leq 1$. Example in notes

One choice to consider is how much variance we want epsilon to have (more for explore)
Knowing how big an acceptance rate to have is 'Optimal Scaling Theory': some say 23.4%

Some maximise expected squared jumping distance $\mathbb{E} [||X_{t+1} - X_t||^2]$

Steps: Test variance for T iterations, measure criterion, change to get better performance on criterion. Called **Adaptive MCMC** (change Kernel K_t)

Using this method actually breaks markov property, as we're using the previous data to make a decision about the new distribution. So need to prove convergence again, though fine if the adapting eventually stops

Can also have Adaptive Gibbs Samplers

Langevin Diffusion $dX_t = \frac{1}{2} \nabla \log \pi(X_t) dt + dB_t$ has a stationary distribution π (if π is smooth)

\Rightarrow **ULA Unadjusted Langevin algorithm** $X^{(t)} = X^{(t-1)} + \frac{\sigma}{2} \nabla \log \pi(X^{(t-1)}) + \sigma W$

MALA Metropolis adjusted Langevin Algorithm $X^* = X^{(t-1)} + \frac{\sigma}{2} \nabla \log \pi(X^{(t-1)}) + \sigma W$

so the Metropolis–Hastings acceptance ratio is

$$\frac{\pi(X^*)}{\pi(X^{(t-1)})} \frac{\mathcal{N}(X^{(t-1)}; X^* + \frac{\sigma}{2} \nabla \log \pi(X^*); \sigma^2)}{\mathcal{N}(X^*; X^{(t-1)} + \frac{\sigma}{2} \nabla \log \pi(X^{(t-1)}); \sigma^2)}. \text{ The sigma } W \text{ pulls more to middle}$$

Sampling and Optimisation ► π is log-concave if $U(x) := -\log \pi(x)$ is convex.

$$q(X^* | X^{(t-1)}) = g(X^*; \varphi(X^{(t-1)})) \quad \text{psi is a deterministic mapping}$$

$$\frac{\pi(X^*) q(X^{(t-1)} | X^*)}{\pi(X^{(t-1)}) q(X^* | X^{(t-1)})} = \frac{\pi(X^*) g(X^{(t-1)}; \varphi(X^*))}{\pi(X^{(t-1)}) g(X^*; \varphi(X^{(t-1)}))}. \quad \text{Can use higher derivatives}$$

Transformations:

Want to sample $\text{supp}(\pi) \subset \mathbb{R}^+$, the posterior distribution of a variance/scale parameter
Any proposed move to R- is a waste of time, so can transform proposals in R to only ones in

\mathbb{R}^+ : Given $X^{(t-1)}$, propose $X^* = \exp(\log X^{(t-1)} + \epsilon)$. Then acceptance =

$$\alpha(X^* | X^{(t-1)}) = \min \left(1, \frac{\pi(X^*)}{\pi(X^{(t-1)})} \frac{q(X^{(t-1)} | X^*)}{q(X^* | X^{(t-1)})} \right) = \min \left(1, \frac{\pi(X^*)}{\pi(X^{(t-1)})} \frac{X^*}{X^{(t-1)}} \right)$$

$$\frac{q(y | x)}{q(x | y)} = \frac{\frac{1}{y\sigma\sqrt{2\pi}} \exp \left[-\frac{(\log y - \log x)^2}{2\sigma^2} \right]}{\frac{1}{x\sigma\sqrt{2\pi}} \exp \left[-\frac{(\log x - \log y)^2}{2\sigma^2} \right]} = \frac{x}{y}.$$

Because

Random Proposals: Acceptance not considering integral

Assume you want to use $q_{\sigma^2}(X^* | X^{(t-1)}) = \mathcal{N}(X^*; X^{(t-1)}, \sigma^2)$
but you don't know how to pick σ^2 . You decide to pick a random $\sigma^{2,*}$ from a distribution $f(\sigma^2)$:

$$\sigma^{2,*} \sim f(\sigma^{2,*}), X^* | \sigma^{2,*} \sim q_{\sigma^{2,*}}(\cdot | X^{(t-1)})$$

$$q(X^* | X^{(t-1)}) = \int q_{\sigma^{2,*}}(X^* | X^{(t-1)}) f(\sigma^{2,*}) d\sigma^{2,*}.$$

Perhaps $q(X^*|X^{(t-1)})$ cannot be evaluated, e.g. the above integral is intractable. Hence the acceptance probability

$$\min \left\{ 1, \frac{\pi(X^*) q(X^{(t-1)}|X^*)}{\pi(X^{(t-1)}) q(X^*|X^{(t-1)})} \right\}$$

cannot be computed.

► Instead you decide to accept your proposal with probability

$$\alpha_t = \min \left\{ 1, \frac{\pi(X^*) q_{\sigma^{2,(t-1)}}(X^{(t-1)}|X^*)}{\pi(X^{(t-1)}) q_{\sigma^{2,*}}(X^*|X^{(t-1)})} \right\}$$

where $\sigma^{2,(t-1)}$ corresponds to parameter of the last accepted proposal.

With probability α_t , set $\sigma^{2,(t)} = \sigma^{2,*}$, $X^{(t)} = X^*$, otherwise $\sigma^{2,(t)} = \sigma^{2,(t-1)}$, $X^{(t)} = X^{(t-1)}$.

Makes MC ($X \times \text{Sigma}^2$)

$$\tilde{\pi}(x, \sigma^2) := \pi(x) f(\sigma^2) \quad q(y, \tau^2|x, \sigma^2) = f(\tau^2) q_{\tau^2}(y|x)$$

$$\frac{\tilde{\pi}(y, \tau^2)}{\tilde{\pi}(x, \sigma^2)} \frac{q(x, \sigma^2|y, \tau^2)}{q(y, \tau^2|x, \sigma^2)} = \frac{\pi(y) f(\tau^2)}{\pi(x) f(\sigma^2)} \frac{f(\sigma^2) q_{\sigma^2}(x|y)}{f(\tau^2) q_{\tau^2}(y|x)} = \frac{\pi(y)}{\pi(x)} \frac{q_{\sigma^2}(x|y)}{q_{\tau^2}(y|x)}$$

Convergence Diagnostics

In general, impossible to see whether MCMC has converged, *but can sometimes know there it hasn't*

Autocorrelogram: Visual diagnostics, shows autocorrelation (correlation of a signal with a delayed copy of itself), almost all MCMC produce positively correlated samples

Geweke's Diagnostics: Uses partial mean of the samples. Determines the burn-in period: the smallest early portion of the chain that passes the diagnostic (i.e. mean of first x%).

Example if the target mean is -2, you ask for the burn in period to be within 0.5 of -2

Gelman Rubin Diagnostics: In order to check for getting stuck at one section. start M chains from various starting points, should be similar

Look at classical sum of squares within group and without group

$$\sum_{m=1}^M \sum_{t=1}^T (X_{m,t} - \bar{X}_{.,.})^2 = \sum_{m=1}^M \sum_{t=1}^T (\bar{X}_{m,.} - \bar{X}_{.,.})^2 \quad \text{inter-group} + \sum_{m=1}^M \sum_{t=1}^T (X_{m,t} - \bar{X}_{m,.})^2 \quad \text{intra-group}$$

Expect inter average inter and intra group value to be the same

$$W = \frac{1}{M} \sum_{m=1}^M \frac{1}{T-1} \sum_{t=1}^T (X_{m,t} - \bar{X}_{m,.})^2 \quad B = \frac{1}{M-1} \sum_{m=1}^M (\bar{X}_{m,.} - \bar{X}_{.,.})^2 \quad V = \left(1 - \frac{1}{T}\right) W + B$$

W and V should both converge to true variance of target distribution, V would be unbiased if starting points were drawn from target, and W usually smaller than V. So plot $\sqrt{V/W}$ and compare with 1, it should converge to 1 as we approach true distribution. If we don't approach it then we've messed up.

Parallel Tempering: run N chains targeting different versions of pi of 'increasing difficulty'.

Introduce **inverse temperatures** $0 < \gamma_1 < \gamma_2 < \dots < \gamma_N = 1$, introduce tempered

distributions π^{γ_n} , and N chains for each. For $y \sim 0$, π^{γ_n} is considered easier to sample because the variations are smaller.

If we want modes of pi, we can just use high temperature chains. If we want to sample from pi, can use high temperature chains to improve mixing of the chain target pi.

Parallel tempering proposes swaps between the different temperature chains.

Joint chain targeting $\pi^{\gamma_1} \otimes \pi^{\gamma_2} \otimes \dots \otimes \pi^{\gamma_N}$. Occasionally perform swap:

Sample indices k_1, k_2 uniformly in $\{1, \dots, N\}$, with acceptance probability

$$\min \left(1, \frac{\pi^{\gamma_{k_1}}(x_{k_2}) \pi^{\gamma_{k_2}}(x_{k_1})}{\pi^{\gamma_{k_1}}(x_{k_1}) \pi^{\gamma_{k_2}}(x_{k_2})} \right).$$

exchange the value of x_{k_1} and x_{k_2} . The swap

moves preserve detailed balance, doesn't change joint target distribution, and the N-th chain still targets π . Doesn't work well in high dimensions when different modes have different shapes

Hamiltonian Monte Carlo (HMC):

$\mu(\mathbf{q}) \propto \exp(-U(\mathbf{q}))$ $U(\mathbf{x})$ is **potential energy**.

$$\exp(-U(\mathbf{q}) - K(\mathbf{p})) = \exp(-U(\mathbf{q})) \exp\left(-\frac{\mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}}{2}\right),$$

decomposition shows **position** $\frac{\mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}}{2}$
q and **momentum p** variables are independent, $K(\mathbf{p})$ is **Kinetic Energy**

$$H(\mathbf{q}, \mathbf{p}) := \frac{\mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}}{2} + U(\mathbf{q}) = K(\mathbf{p}) + U(\mathbf{q}).$$

Hamiltonian Energy need be constant

$$\text{as } \frac{d}{dt} H(\mathbf{q}(t), \mathbf{p}(t)) = \mathbf{p}(t)^T \mathbf{M}^{-1} \frac{d}{dt} \mathbf{p}(t) + \left(\frac{d}{dt} \mathbf{q}(t) \right)^T \nabla U(\mathbf{q}) = 0$$

Rewriting:

$$\begin{aligned} \frac{d}{dt} \mathbf{q} &= +\nabla_{\mathbf{p}} H(\mathbf{q}, \mathbf{p}), \\ \frac{d}{dt} \mathbf{p} &= -\nabla_{\mathbf{q}} H(\mathbf{q}, \mathbf{p}). \end{aligned}$$

t z := $\begin{pmatrix} \mathbf{q} \\ \mathbf{p} \end{pmatrix}$ **J canonical structure matrix** $\mathbf{J} := \begin{pmatrix} \mathbf{0} & \mathbf{I}_d \\ -\mathbf{I}_d & \mathbf{0} \end{pmatrix}$

$$\frac{d}{dt} \mathbf{z} = \mathbf{J} \nabla_{\mathbf{z}} H(\mathbf{z}).$$

Canonical Hamiltonian Equations

Symplectic map: A smooth map $\Psi : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$ s.t.

Jacobian $\nabla_{\mathbf{z}} \Psi(\mathbf{z})$ satisfies that $[\nabla \Psi(\mathbf{z})]^T \mathbf{J}^{-1} \nabla \Psi(\mathbf{z}) = \mathbf{J}^{-1}$

for all \mathbf{z} in \mathbb{R}^{2d} , where $\mathbf{J} = \begin{pmatrix} \mathbf{0} & \mathbf{I}_d \\ -\mathbf{I}_d & \mathbf{0} \end{pmatrix}$ is the canonical structure matrix. This satisfies that $\mathbf{J}^{-1} = \mathbf{J}^T$.

Proposition: Symplectic Maps are volume preserving (so $\det(\mathbf{J}) = 1$)

Proof: infinitesimal cube $[z_1, z_1 + \delta] \times \dots \times [z_{2d}, z_{2d} + \delta]$ maps to thing, volume $\det(\nabla_{\mathbf{z}} \Psi(\mathbf{z})) |\delta^{2d}|$.

By product rule $\det(\nabla_{\mathbf{z}} \Psi(\mathbf{z}))^2 \det(\mathbf{J}^{-1}) = \det(\mathbf{J}^{-1})$, $|\det(\nabla_{\mathbf{z}} \Psi(\mathbf{z}))| = 1$ since $|\det(\mathbf{J}^{-1})| = 1$, so infinitesimal volumes are preserved, same can be obtained for non-infinitesimal by integration

Proposition: Hamiltonian flow is symplectic (and volume preserving)

Let $\Psi_{t,H}(z)$ denote the flow-map of the Hamiltonian dynamics with Hamiltonian $H(z)$, i.e. $\Psi_{t,H}(z(0)) = z(t)$.

We need to understand the behaviour of the Jacobian $\frac{\partial}{\partial z} \Psi_{t,H}(z)$.

We are going to describe the evolution of this Jacobian in time

Let $z(0)$ be an initial point, and $\bar{z}(0) = z(0) + \delta z(0)$ be another nearby initial point.

The differences of two paths at time t can be written as

$$\delta z(t) := \bar{z}(t) - z(t) = \Psi_{t,H}(\bar{z}(0)) - \Psi_{t,H}(z(0)) = \frac{\partial}{\partial z(0)} \Psi_{t,H}(z(0)) \cdot \delta z(0) + o(\delta z(0)).$$

$$\text{Different in } t \quad \frac{d}{dt} \left[\frac{\partial}{\partial z(0)} \Psi_{t,H}(z(0)) \cdot \delta z(0) \right] = \frac{\partial}{\partial z(0)} \left(\frac{d}{dt} \Psi_{t,H}(z(0)) \right) \cdot \delta z(0) = \frac{\partial}{\partial z(0)} (\mathbf{J} \nabla H(z(t))) \cdot \delta z(0)$$

$$= \frac{\partial}{\partial z(0)} (\mathbf{J} \nabla H(\Psi_{t,H}(z(0)))) \cdot \delta z(0) = \mathbf{J} \nabla^2 H(z(t)) \cdot \left[\frac{\partial}{\partial z(0)} \Psi_{t,H}(z(0)) \cdot \delta z(0) \right]$$

Step 2 $F(t) := \frac{\partial}{\partial z(0)} \Psi_{t,H}(z(0))$ = Jacobian at time t , as variation equations hold for every

direction $\frac{d}{dt} F = \mathbf{J} \nabla^2 H(z(t)) F$, $F(0) = I_{2d}$, enough then to verify $F(t)^T \mathbf{J}^{-1} F(t) = I_{2d}$.

Step 3 $F(0) = I_{2d}$, it holds for $t = 0$.

$$\begin{aligned} \frac{d}{dt} (F(t)^T \mathbf{J}^{-1} F(t)) &= (\mathbf{J} \nabla^2 H(z(t)) F(t))^T \mathbf{J}^{-1} F(t) + F(t)^T \mathbf{J}^{-1} \mathbf{J} \nabla^2 H(z(t)) F(t) \\ &= F(t)^T \nabla^2 H(z(t)) \mathbf{J}^T \mathbf{J}^{-1} F(t) + F(t)^T \mathbf{J}^{-1} \mathbf{J} \nabla^2 H(z(t)) F(t) = 0, \text{ since } \mathbf{J}^T \mathbf{J}^{-1} = -\mathbf{J} \mathbf{J}^{-1} = -I_{2d} \text{ and } \mathbf{J}^{-1} \mathbf{J} = I_{2d}. \end{aligned}$$

So Hamiltonian flow is symplectic, so volume preserving by previous proposition

The distribution $\pi(dz) \propto \exp(-H(z)) dz$ is stationary with respect

Proposition to the Hamiltonian flow.

Proof Let the density of π be equal $\pi(z) = \frac{\exp(-H(z))}{C}$ for some normalising constant C .

Let $Z_0 \sim \pi$, and $Z_t = \Psi_{t,H}(Z(0))$. Because $\Psi_{t,H}$ is 1-to-1 mapping

$$Z_t \sim \pi_t, \quad \pi_t(z_t) = \pi(z_0) |\det F(t)|, \quad z_t = \Psi_{t,H}(z_0).$$

By volume and Hamiltonian preserving, $|\det F(t)| = 1$ and $e^{-H(z_0)} = e^{-H(z_t)}$. Then $\pi_t(z_t) = \pi(z_t)$ so $Z_t \sim \pi$.

► Hamilton's equations themselves do not define an ergodic Markov chain, as they keep the Hamiltonian invariant.

Proposition Let $\pi(q) \propto \exp(-U(q))$. Let $T \sim \nu_T$, with positive density on an interval $[0, \tau]$ for $\tau > 0$.

Let K be the Markov kernel for the position variables corresponding to

1. sampling a random momentum p from $\mathcal{N}(0, I_d)$, then

2. running the Hamiltonian dynamics started at (q, p) up to time T sampled from ν_T (independently at each step),

3. and finally discarding the momentum variable.

Suppose that U is continuously differentiable on \mathbb{R}^d , and satisfies that $\sup_q \|\nabla^2 U(q)\| \leq L$, and $\inf_{q \in \mathbb{R}^d} U(q) > -\infty$. Then K is strongly π -irreducible.

Discretizing Hamilton's equations:

Explicit Scheme: $\mathbf{p}(t + \epsilon) = \mathbf{p}(t) - \epsilon \nabla U(\mathbf{q}(t))$ $\mathbf{q}(t + \epsilon) = \mathbf{q}(t) + \epsilon \mathbf{M}^{-1} \mathbf{p}(t)$.

Modified Explicit Scheme: $\mathbf{p}(t + \epsilon) = \mathbf{p}(t) - \epsilon \nabla U(\mathbf{q}(t))$ $\mathbf{q}(t + \epsilon) = \mathbf{q}(t) + \epsilon \mathbf{M}^{-1} \mathbf{p}(t + \epsilon)$.

Leapfrog (Stormer-Verlet Scheme): $\mathbf{p}(t + \epsilon/2) = \mathbf{p}(t) - \frac{\epsilon}{2} \nabla U(\mathbf{q}(t))$

$\mathbf{q}(t + \epsilon) = \mathbf{q}(t) + \epsilon \mathbf{M}^{-1} \mathbf{p}(t + \epsilon/2)$. $\mathbf{p}(t + \epsilon) = \mathbf{p}(t + \epsilon/2) - \frac{\epsilon}{2} \nabla U(\mathbf{q}(t + \epsilon))$

Note: as it is an approximation we don't really expect Hamiltonian energy to always be exactly the same, so we need to recheck the stable solution. Can try to stop it from running away by rejecting the next step if energy is too different

Proposition: The above 3 schemes are symplectic, and hence volume preserving

- ▶ By the chain rule for Jacobians, if $\Psi(z) = \Psi_1(\Psi_2(z))$, then

Proof: $\nabla_z \Psi(z) = \nabla \Psi_1(\Psi_2(z)) \cdot \nabla \Psi_2(z)$.

▶ Assuming Ψ_1 and Ψ_2 are symplectic, then $[\nabla \Psi(z)]^T \mathbf{J}^{-1} \nabla_z \Psi(z) =$

$$(\nabla \Psi_2(z))^T (\nabla \Psi_1(\Psi_2(z)))^T \mathbf{J}^{-1} \nabla \Psi_1(\Psi_2(z)) \cdot \nabla \Psi_2(z) = \mathbf{J}^{-1},$$

By induction this holds for more than two maps, so need to show single step in the schemes

▶ For example, the step $\mathbf{p}(t + \epsilon) = \mathbf{p}(t) - \epsilon \nabla U(\mathbf{q}(t))$ corresponds to the map $\Psi(z) = \begin{pmatrix} \mathbf{q} \\ \mathbf{p} - \epsilon \nabla U(\mathbf{q}) \end{pmatrix}$, which has

$$\text{Jacobian } \nabla \Psi(z) = \begin{pmatrix} \mathbf{I}_d & \mathbf{0} \\ -\epsilon \nabla^2 U(\mathbf{q}) & \mathbf{I}_d \end{pmatrix} \quad \text{▶ By direct calculation, } (\nabla \Psi(z))^T \mathbf{J}^{-1} \nabla \Psi(z)$$

$$= \begin{pmatrix} \mathbf{I}_d & -\epsilon \nabla^2 U(\mathbf{q}) \\ \mathbf{0} & \mathbf{I}_d \end{pmatrix} \begin{pmatrix} \mathbf{0} & -\mathbf{I}_d \\ \mathbf{I}_d & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{I}_d & \mathbf{0} \\ -\epsilon \nabla^2 U(\mathbf{q}) & \mathbf{I}_d \end{pmatrix} = \begin{pmatrix} \mathbf{0} & -\mathbf{I}_d \\ \mathbf{I}_d & \mathbf{0} \end{pmatrix} = \mathbf{J}^{-1}. \quad \text{So symplecticity holds}$$

Hamilton Monte Carlo

- ▶ Let $\Psi_\epsilon : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$ denote the Leapfrog map, and for some $L \in \mathbb{N}$, we denote by Ψ_ϵ^L the L times composition of the Leapfrog map.
- ▶ Let $\mathbf{N} : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$ denote the map that negates the momentum, i.e. $\mathbf{N} \begin{pmatrix} \mathbf{q} \\ \mathbf{p} \end{pmatrix} = \begin{pmatrix} \mathbf{q} \\ -\mathbf{p} \end{pmatrix}$.
- ▶ Let \mathbf{P}_R be the Markov kernel that resamples the momentum component \mathbf{p} from a multivariate Gaussian distribution with covariance matrix \mathbf{M} .
 1. resample the momentum component \mathbf{p} , i.e. apply the Markov kernel \mathbf{P}_R .
 2. We propose a new position $\begin{pmatrix} \mathbf{q}^* \\ \mathbf{p}^* \end{pmatrix} = \Psi \begin{pmatrix} \mathbf{q} \\ \mathbf{p} \end{pmatrix} := \mathbf{N} \left(\Psi_\epsilon^L \begin{pmatrix} \mathbf{q} \\ \mathbf{p} \end{pmatrix} \right)$ by applying L Leapfrog steps and then flipping the momentum. This new position is accepted with probability

$$\min [1, \exp(H(\mathbf{q}, \mathbf{p}) - H(\mathbf{q}^*, \mathbf{p}^*))] = \min \left[1, \exp \left(U(\mathbf{q}) - U(\mathbf{q}^*) + \frac{1}{2} \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p} - \frac{1}{2} (\mathbf{p}^*)^T \mathbf{M}^{-1} \mathbf{p}^* \right) \right].$$

Proposition: π is invariant wrt Markov kernel above, and is reversible wrt second step of the Markov Kernel

- ▶ Since \mathbf{q} and \mathbf{p} are independent according to the target distribution π , it is clear that the first step keeps the target

Proof: invariant.

Let \mathbf{P}_2 denote the Markov kernel corresponding to the second step (a combination of a deterministic step and Metropolis-Hastings accept-reject step), then we are going to check that π is reversible with respect to \mathbf{P}_2 .

- Recall that in general a Markov kernel K is reversible with respect to a distribution π on state space \mathbb{X} if for every bounded measurable function $f : \mathbb{X}^2 \rightarrow \mathbb{R}$, we have

$$\int \int f(\mathbf{x}, \mathbf{y}) \pi(d\mathbf{x}) K(\mathbf{x}, d\mathbf{y}) = \int \int f(\mathbf{x}, \mathbf{y}) \pi(d\mathbf{y}) K(\mathbf{y}, d\mathbf{x}). P_2(\mathbf{x}, d\mathbf{y}) \text{ is non-zero only for } \mathbf{y} = \Psi(\mathbf{x}) \text{ and } \mathbf{y} = \mathbf{x}.$$

$$\int \int f(\mathbf{x}, \mathbf{y}) \pi(d\mathbf{x}) P_2(\mathbf{x}, d\mathbf{y}) = \int \int f(\mathbf{x}, \Psi(\mathbf{x})) \min[1, e^{H(\mathbf{x}) - H(\Psi(\mathbf{x}))}] \pi(d\mathbf{x}) + \int \int f(\mathbf{x}, \mathbf{x}) (1 - \min[1, e^{H(\mathbf{x}) - H(\Psi(\mathbf{x}))}]) \pi(d\mathbf{x})$$

- Let $\mathbf{y} = \Psi(\mathbf{x})$, then $\mathbf{x} = \Psi(\mathbf{y}) = \Psi(\Psi(\mathbf{x}))$. As volume preserving:

$$\pi(d\mathbf{y}) = \pi(d\mathbf{x}) \times \underbrace{|\det(\nabla\Psi)|}_{=1} \frac{\exp(-H(\mathbf{y}))}{\exp(-H(\mathbf{x}))} = \pi(d\mathbf{x}) \cdot e^{H(\mathbf{x}) - H(\mathbf{y})}$$

The first part of the previous sum

$$\text{can be written: } \int \int f(\mathbf{x}, \Psi(\mathbf{x})) \min[1, e^{H(\mathbf{x}) - H(\Psi(\mathbf{x}))}] \pi(d\mathbf{x})$$

$$= \int \int f(\Psi(\mathbf{y}), \mathbf{y}) \min[1, e^{H(\Psi(\mathbf{y})) - H(\mathbf{y})}] \pi(d\mathbf{x}) = \int \int f(\Psi(\mathbf{y}), \mathbf{y}) \min[1, e^{H(\Psi(\mathbf{y})) - H(\mathbf{y})}] \cdot e^{H(\mathbf{y}) - H(\Psi(\mathbf{y}))} \pi(d\mathbf{y})$$

$$= \int \int f(\Psi(\mathbf{y}), \mathbf{y}) \min[1, e^{H(\mathbf{y}) - H(\Psi(\mathbf{y}))}] \pi(d\mathbf{y}).$$

The second part satisfies that

$$\int \int f(\mathbf{x}, \mathbf{x}) (1 - \min[1, e^{H(\mathbf{x}) - H(\Psi(\mathbf{x}))}]) \pi(d\mathbf{x}) = \int \int f(\mathbf{y}, \mathbf{y}) (1 - \min[1, e^{H(\mathbf{y}) - H(\Psi(\mathbf{y}))}]) \pi(d\mathbf{y}).$$

$$\text{Combining } \int \int f(\mathbf{x}, \mathbf{y}) \pi(d\mathbf{x}) K(\mathbf{x}, d\mathbf{y}) = \int \int f(\mathbf{x}, \mathbf{y}) \pi(d\mathbf{y}) K(\mathbf{y}, d\mathbf{x})$$

How to choose epsilon and L: new Markov kernel is a mixture of kernels where π_i is stationary, so π_i is still stationary for it. Using random period length means we can avoid issues that sometimes happen because the deterministic nature of Hamiltonian mechanics

Concentration in High Dimensions:

Suppose that $Z \sim N(0, I_d)$ is a d dimensional standard normal random vector. Then the

Euclidean norm of Z satisfies that for every $t \geq 0$, $\Pr(|\|Z\| - \sqrt{d}| \geq t) \leq C \exp\left(-\frac{t^2}{C}\right)$, C absolute constant

This means that with high probability, $\|Z\| = \sqrt{d} + O(1)$, so most of the probability is concentrated around sphere of radius $d^{0.5}$. Similar examples for other functions (other than Gaussian)

- In general, if the Hessian of the target potential satisfies that $\mu I_d \preceq \nabla^2 U(\mathbf{x}) \preceq L I_d$ for some $0 < \mu < L < \infty$ (strongly convex and smooth potential), and we let $H_{\min} := \inf_z H(z)$, then it is possible to show that

$$\Pr\left(\left|\sqrt{H(z) - H_{\min}} - \mathbb{E}\sqrt{H(z) - H_{\min}}\right| \geq t\right) \leq C \exp\left(-\frac{t^2}{C}\right)$$

Hamiltonian is close to constant high probability density area, and HMC is very efficient in exploring this potentially complicated set automatically.

Optimal Scaling:

Random Walk Metropolis $\mathbf{q}^* = \mathbf{q} + \sigma W$, **MALA** $\mathbf{q}^* = \mathbf{q} + \sigma W + \frac{\sigma^2}{2} \nabla \log \pi(\mathbf{q})$

HMC: $\mathbf{p}(t + \epsilon/2) = \mathbf{p}(t) - \frac{\epsilon}{2} \nabla U(\mathbf{q}(t))$ $\mathbf{q}^* = \mathbf{q} + \epsilon \mathbf{M}^{-1} \mathbf{p}(t + \epsilon/2)$.

$$\mathbf{p}(t + \epsilon) = \mathbf{p}(t + \epsilon/2) - \frac{\epsilon}{2} \nabla U(\mathbf{q}(t + \epsilon)) \quad \text{If } \mathbf{M} = \mathbf{I} \text{ and } \epsilon = \sigma \text{ then one step HMC reduces to MALA.}$$

Maximizing the expected squared jumping distance (ESJD): $\mathbb{E}[\|X^{(t+1)} - X^{(t)}\|^2]$

E.g. Squared jumping distance

$$\|X^{(t+1)} - X^{(t)}\|^2 = \begin{cases} \|\mathbf{q}^* - \mathbf{q}\|^2, & \text{w.p. } \left(1 \wedge \frac{\pi(\mathbf{q}^*) p(\mathbf{q}|\mathbf{q}^*)}{\pi(\mathbf{q}) p(\mathbf{q}^*|\mathbf{q})}\right) \\ 0, & \text{otherwise.} \end{cases}$$

- ▶ RWM: $\mathcal{O}(d^{-1})$
- ▶ MALA: $\mathcal{O}(d^{-1/3})$
- ▶ HMC: $\mathcal{O}(d^{-1/4})$

Dimension dependence on step size :

- ▶ RWM: 0.234 (see Roberts, et. al., 1997)
- ▶ MALA: 0.574 (see Roberts and Rosenthal, 1998)
- ▶ HMC: 0.651 (see Beskos, et. al., 2013)

Optimal Acceptance Rates:

Hamilton Monte Carlo basically non examinable

Hidden Markov Models:

Time Series: models st $Y_t | y_{1:t-1} \sim p(dy_t | y_{1:t-1}, \theta)$, def $\mathbf{y}_{k:l} = (y_k, \dots, y_l)$

$$\forall \theta \in \Theta \quad p(y_{1:t} | \theta) = p(y_1 | \theta) \prod_{k=2}^t p(y_k | y_{1:k-1}, \theta).$$

Likelihood:

can put prior on theta and consider sampling from posterior

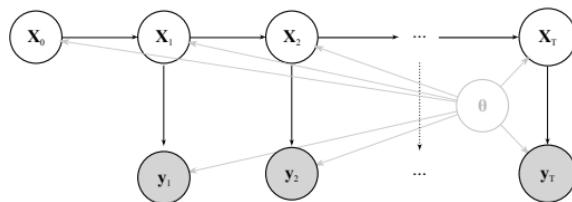
HMMs assume there is a hidden/unobserved X valued Markov Process X_t , with

$X_0 | \theta \sim p(dx_0 | \theta)$, it is time homogenous so that $X_t | x_{t-1}, \theta \sim p(dx_t | x_{t-1}, \theta)$, so

$$p(x_{0:t} | \theta) = p(x_0 | \theta) \prod_{k=1}^t p(x_k | x_{1:k-1}, \theta) = p(x_0 | \theta) \prod_{k=1}^t p(x_k | x_{k-1}, \theta)$$

Observations are dependent on our unobserved Markov process X_t , and conditional

$Y_t | x_t, \theta \sim p(dy_t | x_t, \theta)$ i.e. the distribution of Y_t is independent of $(X_k)_{k \neq t}$ conditional upon $X_t = x_t$.



With a prior $p(\theta)$ on the parameters, the joint density is:

$$\Rightarrow p(y_{1:T} | x_{1:T}, \theta) = \prod_{k=1}^T p(y_k | x_k, \theta) \quad (\text{cond. independent}) \quad p(x_{0:T}, y_{1:T}, \theta) = p(\theta) p(x_0 | \theta) \prod_{t=1}^T p(x_t | x_{t-1}, \theta) \prod_{t=1}^T p(y_t | x_t, \theta).$$

Major goal 1: $p(x_{0:T}, y_{1:T}, \theta)$ or $p(x_t | y_{1:t}, \theta)$ (the **filtering distribution**)

Major goal 2: do **parameter estimation**

Linear Gaussian Models: $\mathbb{X} = \mathbb{R}^{d_x}$ and $\mathbb{Y} = \mathbb{R}^{d_y}$, $X_t = AX_{t-1} + \varepsilon_t$, $\varepsilon \sim \mathcal{N}(0, \Sigma_x)$, need

not be diagonal covariance, $Y_t = CX_t + \eta_t$, $\eta_t \sim \mathcal{N}(0, \Sigma_y)$, can get distributions $X_{1:t}$ and $Y_{1:t}$ by **Kalman Recursions**, and likelihood parameters can be evaluated exactly

$$p(x_{0:t}|y_{1:t}, \theta) = \frac{p(x_{0:t}, y_{1:t} | \theta)}{p(y_{1:t} | \theta)} = \frac{p(x_{0:t} | \theta)p(y_{1:t} | x_{0:t}, \theta)}{p(y_{1:t} | \theta)}$$

General Case:

$$p(x_{0:t} | \theta) = p(x_0 | \theta) \prod_{k=1}^t p(x_k | x_{k-1}, \theta) \quad p(y_{1:t} | x_{0:t}, \theta) = \prod_{k=1}^t p(y_k | x_k, \theta)$$

$$p(y_{1:t} | \theta) = \int_{\mathbb{X}^t} p(x_{0:t}, y_{1:t} | \theta) dx_{0:t}.$$

Marginal Likelihood highD, hard to get

Proposition: The posterior $p(x_{0:t}|y_{1:t}, \theta)$ satisfies

$$p(x_{0:t}|y_{1:t}, \theta) = p(x_{0:t-1}|y_{1:t-1}, \theta) \frac{p(x_t|x_{t-1}, \theta)p(y_t|x_t, \theta)}{p(y_t|y_{1:t-1}, \theta)} \quad (\text{where dropping theta})$$

$$p(y_t|y_{1:t-1}) = \int p(x_{0:t-1}|y_{1:t-1}) p(x_t|x_{t-1}) p(y_t|x_t) dx_{0:t}.$$

$$p(x_{0:t}, y_{1:t}) = p(x_{0:t-1}, y_{1:t-1}) p(x_t|x_{t-1}) p(y_t|x_t)$$

Proof $p(y_{1:t}) = p(y_{1:t-1}) p(y_t|y_{1:t-1})$, and so

$$p(x_{0:t}|y_{1:t}) = \frac{p(x_{0:t}, y_{1:t})}{p(y_{1:t})} = \underbrace{\frac{p(x_{0:t-1}, y_{1:t-1})}{p(y_{1:t-1})}}_{p(x_{0:t-1}|y_{1:t-1})} \underbrace{\frac{p(x_t|x_{t-1})p(y_t|x_t)}{p(y_t|y_{1:t-1})}}_{\text{and the expression for } p(y_t|y_{1:t-1}) \text{ follows}}$$

In general, this means filtering problem is intractable

$$\int \varphi(x_t) p(x_t | y_{1:t}, \theta) dx_t = \int \varphi(x_t) p(x_{0:t} | y_{1:t}, \theta) dx_{0:t} = \frac{\int \varphi(x_t) p(x_0 | \theta) \prod_{k=1}^t p(x_k | x_{k-1}, \theta) \prod_{k=1}^t p(y_k | x_k, \theta) dx_{0:t}}{p(y_{1:t} | \theta)}$$

Numerator is ($t \times \dim(\mathbb{X})$) high dimensional integral, denominator is the marginal likelihood which is also intractable. This **double intractability** means we can't even compute it pointwise for metropolis hastings

$$p(y_{1:t} | \theta) = \int p(x_{0:t}, y_{1:t} | \theta) dx_{0:t} = \int p(x_0 | \theta) \prod_{k=1}^t p(x_k | x_{k-1}, \theta) \prod_{k=1}^t p(y_k | x_k, \theta) dx_{0:t}$$

Historical approach consists of Gibbs Sampling joint space of $(\theta, X_{0:t})$, alternative

between sample from $\theta | x_{0:t}, y_{1:t}$, with conditional distribution

$$p(\theta | x_{0:t}, y_{1:t}) \propto p(\theta) p(x_{0:t}, y_{1:t} | \theta) = p(\theta) p(x_0 | \theta) \prod_{k=1}^t p(x_k | x_{k-1}, \theta) \prod_{k=1}^t p(y_k | x_k, \theta)$$

And sampling from $p(x_{0:t} | y_{1:t}, \theta)$ by iteratively sampling x_k given

x_{k-1}, y_k, x_{k+1} and θ .

$$p(x_k | x_{-k}, y_{1:t}, \theta) = p(x_k | x_{k-1}, y_k, x_{k+1}, \theta) \propto p(x_k | x_{k-1}, \theta) p(y_k, x_{k+1} | x_k, \theta)$$

$$= p(x_k | x_{k-1}, \theta) p(x_{k+1} | x_k, \theta) p(y_k | x_k, \theta)$$

Sometimes we can evaluate this pointwise, and then use MH to update components. As the $X_{0:t}$ are highly correlated Gibbs will fail

Objects of interest:

- predictive checking through $\mathbb{P}(Y_{t+1} \leq y_{t+1} | y_{1:t})$,
- prediction under parameter uncertainty through $p(x_{t+1} | y_{1:t})$,
- sequential model comparison through $p(y_{1:t})$.

$$p(x_{t+1} | y_{1:t}) = \int_{\Theta} \int_{\mathcal{X}^{t+1}} \underbrace{p(x_{t+1} | x_t, \theta)}_{\text{transition}} \underbrace{p(d\theta, dx_{0:t} | y_{1:t})}_{\text{joint posterior}}$$

$$p(y_{1:t}) = \int_{\Theta} \int_{\mathcal{X}^{t+1}} p(d\theta, dx_{0:t}, y_{1:t})$$

Importance Sampling: Know $I = \int \varphi(X) d\pi(X)$, approximate

$$\hat{\pi}^{NIS} = \sum_{i=1}^N W^i \varphi(X^i), W^i = \frac{w(X^i)}{\sum_{i=1}^N w(X^i)}$$

and approximate

$$\pi(x) \approx \pi^N(x) = \sum_{i=1}^N W^i \delta_{X^i}(x)$$

Where X^i are sampled from q and $w(x) = \gamma(x)/q(x)$. , know q isn't write density, just a proposal we can draw from

Sequential Importance Sampling: Key observation that **we can write w(X) as recursion**

$$w(X_{1:t}) = \frac{\gamma(X_{1:t})}{q(X_{1:t})} = \frac{\gamma(X_t | X_{1:t-1}) \gamma(X_{1:t-1})}{q(X_t | X_{1:t-1}) q(X_{1:t-1})} \rightarrow \text{Call } w_t = w(X_{1:t}) \text{ and } \omega_t = \frac{\gamma(X_t | X_{1:t-1})}{q(X_t | X_{1:t-1})}.$$

Then $w_t = w_{t-1} \times \omega_t$

Need 2 additional assumptions: That q, as well as pi, is markovian. And that pi is the

$$q_1(x_1), q_{t|t-1}(X_t | X_{t-1}), \text{ then } \omega_t = \frac{\gamma(X_t | X_{t-1})}{q_{t|t-1}(X_t | X_{t-1})}$$

probability of interest' Moreover, we focus on $\pi(X_{1:t}) = p(X_{1:t} | Y_{1:t})$

Sampling Algorithm for SIS:

At time t = 1: • Sample $X_1^i \sim q_1(\cdot)$. Compute $w_1^i = \frac{\mu(X_1^i) g(y_1 | X_1^i)}{q_1(X_1^i)}$.

At time t > 1: • Sample $X_t^i \sim q_{t|t-1}(\cdot | X_{t-1}^i)$. compute $w_t^i = w_{t-1}^i \times \omega_t^i$

$$= w_{t-1}^i \times \frac{f(X_t^i | X_{t-1}^i) g(y_t | X_t^i)}{q_{t|t-1}(X_t^i | X_{t-1}^i)}.$$

New X_t^i is from Markov chain of q, w_t^i is , the top is out gamma

Outcome: get sequence $X^i = X_{1:t}^i$ and weights $w^i = w_t^i = w(X_{1:t}^i)$ so

$$\pi^N(x) = \sum_{i=1}^N W^i \times \delta_{X^i}(x), \quad W^i = \frac{w^i}{\sum_{j=1}^N w^j}$$

is an approximation of π^N . For example,

$$\begin{aligned}
p(y_t|y_{1:t-1}) &= \int f(x_t|x_{t-1})g(y_t|x_t)p(x_{1:t-1}|y_{1:t-1})dx_{1:t-1}dx_t \\
&= \int \omega_t(x)p(x_{1:t-1}|y_{1:t-1})q_{t|t-1}(x_t|x_{t-1})dx_{1:t-1}dx_t \\
&\approx p^N(y_t|y_{1:t-1}) := \sum_{i=1}^N \omega_t(x^i) \times W_{t-1}^i
\end{aligned}$$

and $p(y_{1:t}) \approx p^N(y_{1:t}) := p^N(y_1) \prod_{k=2}^t p^N(y_k|y_{1:k-1})$, see this by integrating $p(x_{1:t}, y_{1:t})d(x_{1:t})$, using the recursion formula given before

$$q_1(x_1) = \mu(x_1),$$

Prior Proposal: $q_{t|t-1}(x_t|x_{t-1}) = f(x_t|x_{t-1}) \Rightarrow \omega(x_{t-1}, x_t) = g(y_t|x_t)$.

Proposal blindly propagates without considering y_i . We can implement SIS as soon as we can sample from the hidden process $(X_t)_{t \geq 1}$ and evaluate $g(y|x)$ pointwise.

Optimal Proposals: Choose proposal that minimises variance of $(\omega_t^i)_{i=1}^N$

$$q_{t|t-1}^{\text{opt}}(x_t|x_{t-1}) = p(x_t|x_{t-1}, y_t) = \frac{f(x_t|x_{t-1})g(y_t|x_t)}{p(y_t|x_{t-1})} \quad \text{uses observation } y_t \text{ to guide propagation of}$$

$$x_t \Rightarrow \omega_t^{\text{opt}}(x_{t-1}, x_t) = p(y_t|x_{t-1}), \text{ indep of } x_t$$

(proof of optimality in slides) problem: Can't always implement it, e.g. get omega with no closed-form formulae. *Solve via gaussian approximation*

Curse of Dimensionality: SIS as a special type of IS is still cursed by dimensionality

Resampling: motivation: propagating particles with zero weight wastes computational resources,

so at time t , select particle with high weights and remove particles with low weights, spending computational budget on most promising parts

Obtain equally weighted sample (N^{-1}, \bar{X}^i) from a weighted sample (w^i, \tilde{X}^i) . resample on

$$\pi^N(x) = \left(\sum w^j \right)^{-1} \sum w^i \delta_{\tilde{X}^i}(x) \quad \text{output}$$

$$\bar{\pi}^N(x) = N^{-1} \sum \delta_{\bar{X}^i}(x).$$

Multinomial resampling:

Draw ancestry vector

$$A^{1:N} = (A^1, \dots, A^N) \in \{1, \dots, N\}^N \text{ independently from a categorical distribution} \quad A^{1:N} \stackrel{\text{i.i.d}}{\sim} \text{Cat}(w^1, \dots, w^N),$$

Define \bar{X}^i to be X^{A^i} for all $i \in \{1, \dots, N\}$. X^{A^i} is said to be the "parent" or "ancestor" of \bar{X}^i .

$$\text{Return } \bar{X} = (\bar{X}^1, \dots, \bar{X}^N)$$

Forwards version: Draw 'offspring vector' $O^{1:N} = (O^1, \dots, O^N) \in \{0, \dots, N\}^N$

$$O^{1:N} \sim \text{Multinomial}(N; w^1, \dots, w^N) \text{ s.t. } \forall i \in \{1, \dots, N\} \quad \mathbb{E}[O^i] = N \frac{w^i}{\sum_{j=1}^N w^j} \quad \text{and} \quad \sum_{i=1}^N O^i = N.$$

Each particle X^i is replicated O^i times, to create sample \bar{X} :

$$\text{Return } \bar{X} = (\bar{X}^1, \dots, \bar{X}^N)$$

$$\mathbb{E}[O^i] = N \frac{w^i}{\sum_{j=1}^N w^j}, \text{ or } \mathbb{P}[A^i = k] = \frac{w^k}{\sum_{j=1}^N w^j}, \text{ sometimes}$$

Assumed Requirements

$$\frac{1}{N} \sum_{k=1}^N \varphi(\bar{x}^k) = \frac{1}{N} \sum_{k=1}^N O^k \varphi(x^k) \quad \text{has expectation} \quad \sum_{k=1}^N \frac{w^k}{\sum_{j=1}^N w^j} \varphi(x^k)$$

called unbiasedness as \Rightarrow

Therefore, $\bar{\pi}^N$, on average, behaves as the original random measure $\pi^N(x) = \sum W^i \times \delta_{X^i}(x)$. It has more variance!

Sequential Monte Carlo Algorithm

- At time $t = 1$ • Sample $X_1^i \sim q_1(\cdot)$ • Compute the weights

$$w_1^i = \frac{\mu(X_1^i) g(y_1 | X_1^i)}{q_1(X_1^i)}.$$

- At time $t \geq 2$ • Resample $(w_{t-1}^i, X_{1:t-1}^i) \rightarrow (N^{-1}, \bar{X}_{1:t-1}^i)$
- Sample $X_t^i \sim q_{t|t-1}(\cdot | \bar{X}_{t-1}^i)$, $X_{1:t}^i := (\bar{X}_{1:t-1}^i, X_t^i)$ • Compute the weights

$$w_t^i = \omega_t^i = \frac{f(X_t^i | X_{t-1}^i) g(y_t | X_t^i)}{q_{t|t-1}(X_t^i | X_{t-1}^i)}.$$

The resampling breaks the recursion we had in SIS

The cost: **Path Degeneracy**: Sometimes one ancestor leads to more than one offspring, *decreasing diversity*.

Particle approximation of filtering $p(x_t | y_{1:t}, \theta)$: $\frac{1}{\sum_{j=1}^N w_j^i} \sum_{i=1}^N w_t^i \delta_{X_t^i}(dx_t)$.

After resampling $\frac{1}{N} \sum_{i=1}^N \delta_{\bar{X}_t^i}(dx_t)$. Particle approximation of path filtering $P(x_{1:t} | y_{1:t}, \theta)$:

$$\frac{1}{\sum_{j=1}^N w_j^i} \sum_{i=1}^N w_t^i \delta_{X_{1:t}^i}(dx_{1:t}), \quad \bullet \text{ Particle filters approximate well } p(x_t | y_{1:t}) \text{ but not } p(x_s | y_{1:t}) \text{ for } s \ll t.$$

Usually we only care about $(pxt|y1:t)$, so not too bad

Likelihood Estimation: At time 1 (using sample approximation)

$$p^N(y_1) = \frac{1}{N} \sum_{i=1}^N \omega_1^i \xrightarrow[N \rightarrow \infty]{\text{a.s.}} \int \frac{\mu(x_1) g(y_1 | x_1)}{q_1(x_1)} q_1(x_1) dx_1 = p(y_1).$$

$$p^N(y_t | y_{1:t-1}) = \frac{1}{N} \sum_{i=1}^N \omega_t^i$$

At time t (using after sampling approximation)

$$\xrightarrow[N \rightarrow \infty]{\text{a.s.}} \int \omega(x_{t-1}, x_t) q_{t|t-1}(x_t | x_{t-1}) p(x_{t-1} | y_{1:t-1}) dx_{t-1:t} = p(y_t | y_{1:t-1})$$

$$\text{where } \omega(x_{t-1}, x_t) = (f(x_t | x_{t-1}) g(y_t | x_t)) / (q_{t|t-1}(x_t | x_{t-1}))$$

$$p^N(y_{1:t}) = p^N(y_1) \prod_{s=2}^t p^N(y_s | y_{1:s-1}) = \prod_{s=1}^t \frac{1}{N} \sum_{i=1}^N \omega_s^i \xrightarrow[N \rightarrow \infty]{\text{a.s.}} p(y_{1:t}).$$

Leads to estimator

Surprisingly unbiased $\mathbb{E}[p^N(y_{1:t})] = p(y_{1:t})$, whereas for $t >= 2$

(i.e all of the parts are biased but the whole thing isn't)

Theoretical Results: Consistency as $N \rightarrow \infty$ is simple to prove, as each step (propagation, weighting, resampling) is itself consistent

Convergence results include Central Limit Theorems and non-asymptotic results.

Consider $I(\varphi_t) = \int \varphi_t(x_{1:t}) p(x_{1:t} | y_{1:t}) dx_{1:t}$.

• L_p -bound on the path space: $\mathbb{E}[\|I^N(\varphi_t) - I(\varphi_t)\|^p]^{1/p} \leq \frac{B(t)c(p) \|\varphi_t\|_\infty}{\sqrt{N}}$,

$$\sqrt{N}(I^N(\varphi_t) - I(\varphi_t)) \xrightarrow[N \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \sigma_t^2)$$

As expected, $B(t)$ and σ_t^2 typically grow exponentially fast with t . This is the path degeneracy problem.

Consider instead $I(\varphi_t) = \int \varphi_t(x_t) p(x_t | y_{1:t}) dx_t$. L_p -bound:

$$\mathbb{E}[\|I^N(\varphi_t) - I(\varphi_t)\|^p]^{1/p} \leq \frac{B_1 c(p) \|\varphi_t\|_\infty}{\sqrt{N}} \sqrt{N}(I^N(\varphi_t) - I(\varphi_t)) \xrightarrow[N \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \sigma_t^2)$$

For filtering estimates, error is independent of time t : $\sigma_t^2 < \sigma_{\max}^2$ for all t , and B_1 independent of t .

$$p^N(y_{1:t}) = \prod_{s=1}^t \frac{1}{N} \sum_{i=1}^N w_s^i.$$

Consider estimator of marginal likelihood

$$\text{Unbiasedness } \mathbb{E}[p^N(y_{1:t})] = p(y_{1:t}). \text{ Non-asymptotic relative variance } \mathbb{E}\left(\left(\frac{p^N(y_{1:t})}{p(y_{1:t})} - 1\right)^2\right) \leq \frac{B_3 t}{N}.$$

Choose $N = \mathcal{O}(t)$ to control the relative variance.

Curse of dimensionality: Still cursed by dimensionality

Complexity: Propagating and weighting the particles is $\mathcal{O}(N)$.

Multinomial resampling is $\mathcal{O}(N)$ if the uniforms are generated in sorted order. The memory cost is $\mathcal{O}(N)$ if only the latest particles are stored.

The memory cost is at most $\mathcal{O}(Nt)$ if the paths are stored

Sequential Monte Carlo in Generic Models:

$$\pi_T(\theta) \propto \pi_0(\theta) \prod_{s=1}^T p(y_s | \theta).$$

Want to sample from

, not imposing any structure on theta

here Idea: sample from π_0 and apply SIS-like recursion

$$\text{If } w_t^i = \frac{\pi_t(\theta^i)}{\pi_0(\theta^i)}, \text{ then } w_t^i = \frac{\pi_{t-1}(\theta^i)}{\pi_0(\theta^i)} \frac{\pi_t(\theta^i)}{\pi_{t-1}(\theta^i)} = w_{t-1}^i \omega_t^i, \quad \text{with } \omega_t^i = \frac{\pi_t(\theta^i)}{\pi_{t-1}(\theta^i)}.$$

Note all these methods work even when we haven't normalised weights. Incremental

$$\text{weights satisfy } \omega^i \propto \frac{\pi_t(\theta)}{\pi_{t-1}(\mu)} \propto p(y_t | \theta)$$

Algorithm: • Here X corresponds exactly to θ . No time structure.

- At time $t = 0$ • Sample $X^i \sim \pi_0(\cdot)$ for $i \in \{1, \dots, N\}$ $w_0^i = \frac{1}{N}$
- At time $t \geq 1$ $w_t^i = w_{t-1}^i \times \omega_t^i = w_{t-1}^i \times \frac{\pi_t(X^i)}{\pi_{t-1}(X^i)}$ At all times, $(w_t^i, X^i)_{i=1}^N$ approximates π_t .

Not sampling again at each step; so exacerbate problems of path degeneracy, ending up with very few unique values

SIS with resampling for posterior inference:

- At time $t = 0$ Sample $X_0^i \sim \pi_0(\cdot)$ for $i \in \{1, \dots, N\}$ $w_0^i = \frac{1}{N}$
- At time $t \geq 1$ Resample $(w_{t-1}^i, X_{t-1}^i) \rightarrow (N^{-1}, \bar{X}_{t-1}^i)$ Define $X_t^i = \bar{X}_{t-1}^i$ $w_t^i = \omega_t^i = \frac{\pi_t(X_t^i)}{\pi_{t-1}(X_t^i)}$

Suppose we want to use MH for $p(\text{theta}|y1:t)$, need to compute acceptance ratio which needs $p(y1:t, \text{theta})$, but that's hard to compute. **BUT** can get unbiased estimators of $p(y1:t, \text{theta})$, $P^N(y1:t, \text{theta})$ based on Sequential Monte Carlo (SMC)

Pseudo Marginal Metropolis Hastings: Uses estimates of $\tilde{\pi}(x)$

- Starting with $X^{(1)}$, and $Z^{(1)} \geq 0$ such that $\mathbb{E}(Z^{(1)}) = \tilde{\pi}(X^{(1)})$, for $t = 2, 3, \dots$
- 1. Sample $X' \sim q(\cdot | X^{(t-1)})$ 2. Estimate $\tilde{\pi}(X')$ by Z' , such that $\mathbb{E}(Z') = \tilde{\pi}(X')$.
- 3. Compute $\alpha(X' | X^{(t-1)}) = \min\left(1, \frac{Z' q(X^{(t-1)} | X')}{Z^{(t-1)} q(X' | X^{(t-1)})}\right)$
- 4. Sample $U \sim U_{[0,1]}$. If $U \leq \alpha(X' | X^{(t-1)})$, set $(X^{(t)}, Z^{(t)}) = (X', Z')$, otherwise set $(X^{(t)}, Z^{(t)}) = (X^{(t-1)}, Z^{(t-1)})$.

Exact limiting law

- For any x , denote by Z_x a non-negative unbiased estimator of $\tilde{\pi}(x)$, with distribution $g(\cdot | x) \equiv g_x$.
- Thus the generated chain $(X^{(t)})_{t \geq 1}$ goes to $\approx \pi$.
- If $\mathbb{V}_{g(\cdot|x)}(Z_x/\tilde{\pi}(x)) \ll 1$, then the algorithm \approx original Metropolis–Hastings.
- In fact, the limiting law of $(X^{(t)})_{t \geq 0}$ is exactly π ...

For general: introduce extended target distribution $\bar{\pi}(x, z) \propto z \times g_x(z)$

Introduce proposal $\bar{q}(x', z' | x, z) = q(x' | x)g_{x'}(z')$. \Rightarrow MH acceptance ratio is

$$\min\left(1, \frac{\bar{\pi}(x', z')}{\bar{\pi}(x, z)} \frac{\bar{q}(x, z | x', z')}{\bar{q}(x', z' | x, z)}\right) = \min\left(1, \frac{z' q(x | x')}{z q(x' | x)}\right)$$

This is the algorithm above

Is standard MH targeting $\bar{\pi}(x, z)$. Distribution of X if X, Z following $\bar{\pi}(x, z)$.

Integrating X out $\bar{\pi}(x) \propto \int z g_x(z) dz = \mathbb{E}_{g_x}[Z_x] = \tilde{\pi}(x)$, so marginal is π

- Thus if the Markov chain $(X^{(t)}, Z^{(t)})_{t \geq 0}$ converges to $\bar{\pi}$, then the first component $(X^{(t)})$ converges to the first marginal of $\bar{\pi}$, which is π .

- Therefore pseudo-marginal Metropolis–Hastings is *exact*.

Metropolis Hastings ratio for PMMH algorithm

- Recall that we have

$$p(\theta, x_{1:t} | y_{1:t}) = \frac{1}{p(y_{1:t})} p(\theta) p(x_{1:t} | \theta) p(y_{1:t} | x_{1:t}, \theta) \text{ where}$$

$$p(y_{1:t}) = \int p(\theta) p(y_{1:t} | \theta) d\theta \text{ and } p(y_{1:t} | \theta) = \int p(y_{1:t} | x_{1:t}, \theta) p(x_{1:t} | \theta) dx_{1:t}$$

- Therefore, $p(\theta | y_{1:t}) = \frac{p(\theta) p(y_{1:t} | \theta)}{p(y_{1:t})}$ and $p(y_{1:t} | \theta) \approx p^N(y_{1:t} | \theta)$
- Here $z = p(y_{1:t} | \theta)$ is approximated by $p^N(y_{1:t} | \theta)$ where the randomness is induced by the sampled paths $X_{1:t}$ (targeting the posterior).
- Parameters are drawn from the Kernel $q(\theta' | \theta)$ and $X'_{1:t}$ $r = \frac{p(\theta') p^N_{x'}(y_{1:t} | \theta') q(\theta | \theta')}{p(\theta) p^N_x(y_{1:t} | \theta) q(\theta' | \theta)}$

Gibbs Samplers for HMM: not great

- Sampling from $p(x_{0:t} | y_{1:t}, \theta)$ can be done by iteratively sampling x_k given x_{k-1}, y_k, x_{k+1} and θ .

$$\begin{aligned} p(x_k | x_{-k}, y_{1:t}, \theta) &= p(x_k | x_{k-1}, y_k, x_{k+1}, \theta) \propto p(x_k | x_{k-1}, \theta) p(y_k, x_{k+1} | x_k, \theta) \\ &= p(x_k | x_{k-1}, \theta) p(x_{k+1} | x_k, \theta) p(y_k | x_k, \theta) \end{aligned}$$

- In which case, we could use Metropolis–Hastings to update each component of $X_{0:t}$, given the others.
- By definition, the components of $X_{0:t}$ are highly correlated, thus this Gibbs sampling approach will fail (remember the bivariate normal!).

“Idealised” Gibbs Sampler

- To sample from $p(\theta, x_{1:T} | y_{1:T})$, an MCMC strategy consists of using the following block Gibbs sampler.

At iteration $i = 0$

- Sample $\theta(0), X_{1:T}(0)$ and its ancestral lineage $B_{1:T}(0)$ arbitrarily.

At iteration $i \geq 1$

- Sample $X_{1:T}(i) \sim p_{\theta(i-1)}(x_{1:T} | y_{1:T})$.
- Sample $\theta(i) \sim p(\theta | y_{1:T}, X_{1:T}(i))$
- **Problem:** We do not know how to sample from $p_\theta(x_{1:T} | y_{1:T})$
 - Naive particle approximation where $X_{1:T}(i) \sim \hat{p}(x_{1:T} | y_{1:T}, \theta(i))$ is substituted to $X_{1:T}(i) \sim p(x_{1:T} | y_{1:T}, \theta(i))$ is obviously incorrect.

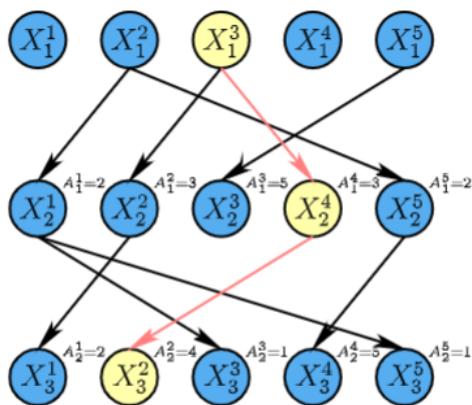
Particle Gibbs Sampler:

At iteration $i = 0$ • Sample $\theta(0)$, $X_{1:T}(0)$ and its ancestral lineage $B_{1:T}(0)$ arbitrarily.

At iteration $i \geq 0$ • Sample $\theta(i) \sim p(\theta|y_{1:T}, X_{1:T}(i-1))$.

- Run a conditional SMC algorithm for $\theta(i)$ consistent with the $X_{1:T}(i-1)$ and its ancestral lineage $B_{1:T}(i-1)$.
- Re-sample the trajectories at T and define uniformly an index in $k \in 1 : N$. Then, pull a trajectory from that index X_T^k and going back it time through its genealogy B_t to define $X_{1:T}(i)$.
- **Proposition.** Assume that the 'ideal' Gibbs sampler chain is ergodic then under very weak assumptions the particle Gibbs sampler chain is ergodic and admits $p(\theta, x_{1:T}|y_{1:T})$ as an invariant distribution for any $N \geq 2$.

Conditional SMC



• In our case we use SMC to approximate $\hat{\pi}(X|\theta)$ and then update θ with a MH move. As we pull a trajectory we go from approximation to 'true' distribution.

- Suppose we have an approximation for the posterior over (θ, X) , $\hat{\pi}(\theta, X)$ based on n particles X^i .
- We can conceive an enlarged space with an index k and define $\hat{\pi}(\theta, X, k)$ as selecting one of the particles.
- If K is random It can be shown that (θ, X^K) distributes according to $\pi(\theta, X)$.