

Algorithms for unifying statistical inference

Ryan Chan

9 Feb 2021



**The
Alan Turing
Institute**

Outline

The fusion problem

Popular algorithms for fusion

The Monte Carlo Fusion algorithm

Constructing a rejection sampler

Simple examples

Possible extensions to Monte Carlo Fusion

Hierarchical Monte Carlo Fusion

Divide-and-Conquer SMC with Fusion

Bayesian Fusion

Fusion Problem

- Target:

$$\pi(\mathbf{x}) \propto \prod_{c=1}^C f_c(\mathbf{x})$$

where each *sub-posterior*, $f_c(\mathbf{x})$, is a density representing one of the C distributed inferences we wish to unify

- Assume we can sample $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$
- Applications:
 - Big Data (by construction)
 - Tempering (by construction)
 - Expert elicitation: combining views of multiple experts
 - Privacy setting

Fusion Problem

- Target:

$$\pi(\mathbf{x}) \propto \prod_{c=1}^C f_c(\mathbf{x})$$

where each *sub-posterior*, $f_c(\mathbf{x})$, is a density representing one of the C distributed inferences we wish to unify

- Assume we can sample $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$
- Applications:
 - Big Data (by construction)
 - Tempering (by construction)
 - Expert elicitation: combining views of multiple experts
 - Privacy setting

Fusion Problem

- Target:

$$\pi(\mathbf{x}) \propto \prod_{c=1}^C f_c(\mathbf{x})$$

where each *sub-posterior*, $f_c(\mathbf{x})$, is a density representing one of the C distributed inferences we wish to unify

- Assume we can sample $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$
- Applications:
 - Big Data (by construction)
 - Tempering (by construction)
 - Expert elicitation: combining views of multiple experts
 - Privacy setting

Fusion for Big Data

- Consider we have data \mathbf{x} with a large number of observations n
- The likelihood $\ell(\mathbf{x} | \theta)$ becomes expensive to calculate
 - This makes MCMC prohibitively slow for big data
- Potential solution:

$$\pi(\theta | \mathbf{x}) \propto \prod_{i=1}^n \ell(\mathbf{x}_i | \theta) \pi(\theta) = \prod_{c=1}^C \ell(\mathbf{x}_c | \theta) \pi(\theta)^{\frac{1}{C}}$$

where \mathbf{x}_c denotes the c -th subset for $c = 1, \dots, C$ and $\pi(\theta) = \prod_{c=1}^C \pi(\theta)^{\frac{1}{C}}$ is the prior

- Advantage: inference on each smaller dataset can be conducted independently and in parallel

Fusion for Big Data

- Consider we have data \mathbf{x} with a large number of observations n
- The **likelihood** $\ell(\mathbf{x} | \theta)$ becomes **expensive** to calculate
 - This makes MCMC prohibitively **slow** for big data
- **Potential solution:**

$$\pi(\theta | \mathbf{x}) \propto \prod_{i=1}^n \ell(\mathbf{x}_i | \theta) \pi(\theta) = \prod_{c=1}^C \ell(\mathbf{x}_c | \theta) \pi(\theta)^{\frac{1}{C}}$$

where \mathbf{x}_c denotes the c -th subset for $c = 1, \dots, C$ and $\pi(\theta) = \prod_{c=1}^C \pi(\theta)^{\frac{1}{C}}$ is the prior

- Advantage: inference on each smaller dataset can be conducted independently and in **parallel**

Fusion for Big Data

- Consider we have data \mathbf{x} with a large number of observations n
- The **likelihood** $\ell(\mathbf{x} | \theta)$ becomes **expensive** to calculate
 - This makes MCMC prohibitively **slow** for big data
- **Potential solution:**

$$\pi(\theta | \mathbf{x}) \propto \prod_{i=1}^n \ell(\mathbf{x}_i | \theta) \pi(\theta) = \prod_{c=1}^C \ell(\mathbf{x}_c | \theta) \pi(\theta)^{\frac{1}{C}}$$

where \mathbf{x}_c denotes the c -th subset for $c = 1, \dots, C$ and $\pi(\theta) = \prod_{c=1}^C \pi(\theta)^{\frac{1}{C}}$ is the prior

- Advantage: inference on each smaller dataset can be conducted independently and in **parallel**

Fusion for Big Data

- Consider we have data \mathbf{x} with a large number of observations n
- The **likelihood** $\ell(\mathbf{x} | \theta)$ becomes **expensive** to calculate
 - This makes MCMC prohibitively **slow** for big data
- **Potential solution:**

$$\pi(\theta | \mathbf{x}) \propto \prod_{i=1}^n \ell(\mathbf{x} | \theta) \pi(\theta) = \prod_{c=1}^C \ell(\mathbf{x}_c | \theta) \pi(\theta)^{\frac{1}{C}}$$

where \mathbf{x}_c denotes the c -th subset for $c = 1, \dots, C$ and $\pi(\theta) = \prod_{c=1}^C \pi(\theta)^{\frac{1}{C}}$ is the prior

- Advantage: inference on each smaller dataset can be conducted independently and in **parallel**

Fusion for Big Data

- Consider we have data \mathbf{x} with a large number of observations n
- The **likelihood** $\ell(\mathbf{x} | \theta)$ becomes **expensive** to calculate
 - This makes MCMC prohibitively **slow** for big data
- **Potential solution:**

$$\pi(\theta | \mathbf{x}) \propto \prod_{i=1}^n \ell(\mathbf{x} | \theta) \pi(\theta) = \prod_{c=1}^C \ell(\mathbf{x}_c | \theta) \pi(\theta)^{\frac{1}{C}}$$

where \mathbf{x}_c denotes the c -th subset for $c = 1, \dots, C$ and $\pi(\theta) = \prod_{c=1}^C \pi(\theta)^{\frac{1}{C}}$ is the prior

- Advantage: inference on each smaller dataset can be conducted independently and in **parallel**

Fusion for Tempering

- Consider the *power-tempered target distribution*
 $\pi_\beta(\mathbf{x}) = [\pi(\mathbf{x})]^\beta$ for $\beta \in (0, 1]$
- MCMC can become computationally expensive to sample from *multi-modal densities* and can get stuck in modes
- Potential solution:

$$\pi(\mathbf{x}) \propto \pi(\mathbf{x})^{\frac{1}{\beta} \cdot \beta} \propto \prod_{c=1}^{\frac{1}{\beta}} \pi_\beta(\mathbf{x})$$

where $\frac{1}{\beta} \in \mathbb{N}$

Fusion for Tempering

- Consider the *power-tempered target distribution*
 $\pi_\beta(\mathbf{x}) = [\pi(\mathbf{x})]^\beta$ for $\beta \in (0, 1]$
- MCMC can become computationally expensive to sample from **multi-modal densities** and can get stuck in modes
- Potential solution:

$$\pi(\mathbf{x}) \propto \pi(\mathbf{x})^{\frac{1}{\beta} \cdot \beta} \propto \prod_{c=1}^{\frac{1}{\beta}} \pi_\beta(\mathbf{x})$$

where $\frac{1}{\beta} \in \mathbb{N}$

Fusion for Tempering

- Consider the *power-tempered target distribution*
 $\pi_\beta(\mathbf{x}) = [\pi(\mathbf{x})]^\beta$ for $\beta \in (0, 1]$
- MCMC can become computationally expensive to sample from **multi-modal densities** and can get stuck in modes
- **Potential solution:**

$$\pi(\mathbf{x}) \propto \pi(\mathbf{x})^{\frac{1}{\beta} \cdot \beta} \propto \prod_{c=1}^{\frac{1}{\beta}} \pi_\beta(\mathbf{x})$$

where $\frac{1}{\beta} \in \mathbb{N}$

Fusion in a privacy setting

- Suppose have C parties that wish to combine their inferences but either:
 - underlying model $f_c(\mathbf{x})$ cannot be shared, or
 - underlying data \mathbf{x}_c cannot be shared
- e.g. healthcare settings
- If we have a method that can combine inferences or samples, would need some mechanism that ensures the privacy of the model or data

Fusion in a privacy setting

- Suppose have C parties that wish to combine their inferences but either:
 - underlying model $f_c(\mathbf{x})$ cannot be shared, or
 - underlying data \mathbf{x}_c cannot be shared
- e.g. healthcare settings
- If we have a method that can combine inferences or samples, would need some mechanism that ensures the privacy of the model or data

Fusion in a privacy setting

- Suppose have C parties that wish to combine their inferences but either:
 - underlying model $f_c(\mathbf{x})$ cannot be shared, or
 - underlying data \mathbf{x}_c cannot be shared
- e.g. healthcare settings
- If we have a method that can combine inferences or samples, would need some mechanism that ensures the privacy of the model or data

Fusion in a privacy setting

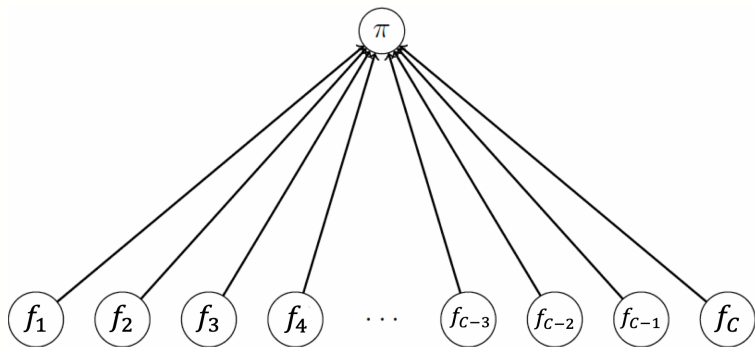
- Suppose have C parties that wish to combine their inferences but either:
 - underlying model $f_c(\mathbf{x})$ cannot be shared, or
 - underlying data \mathbf{x}_c cannot be shared
- e.g. healthcare settings
- If we have a method that can combine inferences or samples, would need some mechanism that ensures the privacy of the model or data

Fusion in a privacy setting

- Suppose have C parties that wish to combine their inferences but either:
 - underlying model $f_c(\mathbf{x})$ cannot be shared, or
 - underlying data \mathbf{x}_c cannot be shared
- e.g. healthcare settings
- If we have a method that can combine inferences or samples, would need some mechanism that ensures the privacy of the model or data

Fork-and-join fusion

The **fork-and-join** approach:



Current Fork-and-Join Methods

- Several fork-and-join methods have been developed. For instance
 - Gaussian approximations to sub-posteriors [Neiswanger et al., 2013]
 - Kernel density averaging [Neiswanger et al., 2013]
 - Consensus Monte Carlo [Scott et al., 2016]

Current Fork-and-Join Methods

- Several fork-and-join methods have been developed. For instance
 - Gaussian approximations to sub-posteriors [Neiswanger et al., 2013]
 - Kernel density averaging [Neiswanger et al., 2013]
 - Consensus Monte Carlo [Scott et al., 2016]

Current Fork-and-Join Methods

- Several fork-and-join methods have been developed. For instance
 - Gaussian approximations to sub-posteriors [Neiswanger et al., 2013]
 - Kernel density averaging [Neiswanger et al., 2013]
 - Consensus Monte Carlo [Scott et al., 2016]

Current Fork-and-Join Methods

- Several fork-and-join methods have been developed. For instance
 - Gaussian approximations to sub-posteriors [Neiswanger et al., 2013]
 - Kernel density averaging [Neiswanger et al., 2013]
 - Consensus Monte Carlo [Scott et al., 2016]

Kernel density averaging (KDEMC)

- Apply a **kernel density estimation** to each sub-posterior, $\hat{f}_c(\mathbf{x})$ [Neiswanger et al., 2013]. Then approximate full posterior by

$$\hat{\pi}(\mathbf{x}) = \prod_{c=1}^C \hat{f}_c(\mathbf{x})$$

- If Gaussian kernels are used, $\hat{\pi}(\mathbf{x})$ is a product of Gaussian mixtures with $O(NC)$ components (N samples, C sub-posteriors)
- Neiswanger et al. [2013] suggest sampling from the Gaussian mixture using MCMC
- Can be computationally expensive and inefficient
- Does not scale well with dimension

Kernel density averaging (KDEMC)

- Apply a **kernel density estimation** to each sub-posterior, $\hat{f}_c(\mathbf{x})$ [Neiswanger et al., 2013]. Then approximate full posterior by

$$\hat{\pi}(\mathbf{x}) = \prod_{c=1}^C \hat{f}_c(\mathbf{x})$$

- If Gaussian kernels are used, $\hat{\pi}(\mathbf{x})$ is a product of Gaussian mixtures with $O(NC)$ components (N samples, C sub-posteriors)
- Neiswanger et al. [2013] suggest sampling from the Gaussian mixture using MCMC
- Can be computationally expensive and inefficient
- Does not scale well with dimension

Kernel density averaging (KDEMC)

- Apply a **kernel density estimation** to each sub-posterior, $\hat{f}_c(\mathbf{x})$ [Neiswanger et al., 2013]. Then approximate full posterior by

$$\hat{\pi}(\mathbf{x}) = \prod_{c=1}^C \hat{f}_c(\mathbf{x})$$

- If Gaussian kernels are used, $\hat{\pi}(\mathbf{x})$ is a product of Gaussian mixtures with $O(NC)$ components (N samples, C sub-posteriors)
- Neiswanger et al. [2013] suggest sampling from the Gaussian mixture using MCMC
- Can be computationally expensive and inefficient
- Does not scale well with dimension

Kernel density averaging (KDEMC)

- Apply a **kernel density estimation** to each sub-posterior, $\hat{f}_c(\mathbf{x})$ [Neiswanger et al., 2013]. Then approximate full posterior by

$$\hat{\pi}(\mathbf{x}) = \prod_{c=1}^C \hat{f}_c(\mathbf{x})$$

- If Gaussian kernels are used, $\hat{\pi}(\mathbf{x})$ is a product of Gaussian mixtures with $O(NC)$ components (N samples, C sub-posteriors)
- Neiswanger et al. [2013] suggest sampling from the Gaussian mixture using MCMC
- Can be computationally expensive and inefficient
- Does not scale well with dimension

Kernel density averaging (KDEMC)

- Apply a **kernel density estimation** to each sub-posterior, $\hat{f}_c(\mathbf{x})$ [Neiswanger et al., 2013]. Then approximate full posterior by

$$\hat{\pi}(\mathbf{x}) = \prod_{c=1}^C \hat{f}_c(\mathbf{x})$$

- If Gaussian kernels are used, $\hat{\pi}(\mathbf{x})$ is a product of Gaussian mixtures with $O(NC)$ components (N samples, C sub-posteriors)
- Neiswanger et al. [2013] suggest sampling from the Gaussian mixture using MCMC
- Can be computationally expensive and inefficient
- Does not scale well with dimension

Consensus Monte Carlo

- Approximate the full posterior as a weighted average of the sub-posterior samples [Scott et al., 2016]
- Suppose have MCMC samples $\mathbf{x}_1^{(c)}, \dots, \mathbf{x}_N^{(c)}$ from $f_c(\mathbf{x})$ for $c = 1, \dots, C$. Then approximate full posterior

$$\mathbf{x}_i = \left(\sum_{c=1}^C W_c \right)^{-1} \left(\sum_{c=1}^C W_c \mathbf{x}_i^{(c)} \right)$$

where $W_c \in \mathbb{R}^d$ is a weight matrix for sub-posterior c (typically take $W_c = \hat{\Sigma}_c$)

- Method is exact if sub-posteriors are Gaussian (motivated by Bernstein-von Mises Theorem)
- Very scalable

Consensus Monte Carlo

- Approximate the full posterior as a weighted average of the sub-posterior samples [Scott et al., 2016]
- Suppose have MCMC samples $\mathbf{x}_1^{(c)}, \dots, \mathbf{x}_N^{(c)}$ from $f_c(\mathbf{x})$ for $c = 1, \dots, C$. Then approximate full posterior

$$\mathbf{x}_i = \left(\sum_{c=1}^C W_c \right)^{-1} \left(\sum_{c=1}^C W_c \mathbf{x}_i^{(c)} \right)$$

where $W_c \in \mathbb{R}^d$ is a weight matrix for sub-posterior c (typically take $W_c = \hat{\Sigma}_c$)

- Method is exact if sub-posteriors are Gaussian (motivated by Bernstein-von Mises Theorem)
- Very scalable

Consensus Monte Carlo

- Approximate the full posterior as a weighted average of the sub-posterior samples [Scott et al., 2016]
- Suppose have MCMC samples $\mathbf{x}_1^{(c)}, \dots, \mathbf{x}_N^{(c)}$ from $f_c(\mathbf{x})$ for $c = 1, \dots, C$. Then approximate full posterior

$$\mathbf{x}_i = \left(\sum_{c=1}^C W_c \right)^{-1} \left(\sum_{c=1}^C W_c \mathbf{x}_i^{(c)} \right)$$

where $W_c \in \mathbb{R}^d$ is a weight matrix for sub-posterior c
(typically take $W_c = \hat{\Sigma}_c$)

- Method is exact if sub-posteriors are Gaussian (motivated by Bernstein-von Mises Theorem)
- Very scalable

Consensus Monte Carlo

- Approximate the full posterior as a weighted average of the sub-posterior samples [Scott et al., 2016]
- Suppose have MCMC samples $\mathbf{x}_1^{(c)}, \dots, \mathbf{x}_N^{(c)}$ from $f_c(\mathbf{x})$ for $c = 1, \dots, C$. Then approximate full posterior

$$\mathbf{x}_i = \left(\sum_{c=1}^C W_c \right)^{-1} \left(\sum_{c=1}^C W_c \mathbf{x}_i^{(c)} \right)$$

where $W_c \in \mathbb{R}^d$ is a weight matrix for sub-posterior c
(typically take $W_c = \hat{\Sigma}_c$)

- Method is exact if sub-posteriors are Gaussian (motivated by Bernstein-von Mises Theorem)
- Very scalable

Current Fork-and-Join Methods

- Several fork-and-join methods have been developed. For instance
 - Gaussian approximations to sub-posteriors [Neiswanger et al., 2013]
 - Kernel density averaging [Neiswanger et al., 2013]
 - Consensus Monte Carlo [Scott et al., 2016]
- A primary weakness of these methods is that the recombination is **inexact** in general and involve **approximations**
- However, Monte Carlo Fusion [Dai et al., 2019] is **exact**

Current Fork-and-Join Methods

- Several fork-and-join methods have been developed. For instance
 - Gaussian approximations to sub-posteriors [Neiswanger et al., 2013]
 - Kernel density averaging [Neiswanger et al., 2013]
 - Consensus Monte Carlo [Scott et al., 2016]
- A primary weakness of these methods is that the recombination is **inexact** in general and involve **approximations**
- However, Monte Carlo Fusion [Dai et al., 2019] is **exact**

Current Fork-and-Join Methods

- Several fork-and-join methods have been developed. For instance
 - Gaussian approximations to sub-posteriors [Neiswanger et al., 2013]
 - Kernel density averaging [Neiswanger et al., 2013]
 - Consensus Monte Carlo [Scott et al., 2016]
- A primary weakness of these methods is that the recombination is **inexact** in general and involve **approximations**
- However, Monte Carlo Fusion [Dai et al., 2019] is **exact**

- └ The Monte Carlo Fusion algorithm
 - └ Constructing a rejection sampler

Constructing a rejection sampler - An Extended Target

Proposition

Suppose that $p_c(\mathbf{y} \mid \mathbf{x}^{(c)})$ is the transition density of a *stochastic process with stationary distribution* $f_c^2(\mathbf{x})$. The $(C + 1)d$ -dimensional (fusion) density proportional to the integrable function

$$g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C \left[f_c^2(\mathbf{x}^{(c)}) p_c(\mathbf{y} \mid \mathbf{x}^{(c)}) \cdot \frac{1}{f_c(\mathbf{y})} \right]$$

admits the *marginal density* π for \mathbf{y} .

- **Main idea:** If we can sample from g , then we can obtain a draw from the fusion density ($\mathbf{y} \sim \pi$)

Constructing a rejection sampler - An Extended Target

Proposition

Suppose that $p_c(\mathbf{y} \mid \mathbf{x}^{(c)})$ is the transition density of a *stochastic process with stationary distribution* $f_c^2(\mathbf{x})$. The $(C + 1)d$ -dimensional (fusion) density proportional to the integrable function

$$g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y}) \propto \prod_{c=1}^C \left[f_c^2(\mathbf{x}^{(c)}) p_c(\mathbf{y} \mid \mathbf{x}^{(c)}) \cdot \frac{1}{f_c(\mathbf{y})} \right]$$

admits the *marginal density* π for \mathbf{y} .

- **Main idea:** If we can sample from g , then we can obtain a draw from the fusion density ($\mathbf{y} \sim \pi$)

Rejection Sampling (Double Langevin Approach)

- There are many possible choices for $p_c(\mathbf{y} \mid \mathbf{x})$
- Let $p_c(\mathbf{y} \mid \mathbf{x}) := p_{T,c}(\mathbf{y} \mid \mathbf{x})$, the transition density of the d -dimensional (double) Langevin (DL) diffusion processes $\mathbf{X}_t^{(c)}$ for $c = 1, \dots, C$, from \mathbf{x} to \mathbf{y} for a pre-defined time $T > 0$ given by

$$d\mathbf{X}_t^{(c)} = \Lambda_c \nabla \log f_c \left(\mathbf{x}_t^{(c)} \right) dt + \Lambda_c^{1/2} d\mathbf{W}_t^{(c)}$$

- $\mathbf{W}_t^{(c)}$ is d -dimensional Brownian motion
- Λ_c is the pre-conditioning matrix associated with sub-posterior f_c
- ∇ is the gradient operator over \mathbf{x}
- Has stationary distribution $f_c^2(\mathbf{x})$

Rejection Sampling (Double Langevin Approach)

- There are many possible choices for $p_c(\mathbf{y} \mid \mathbf{x})$
- Let $p_c(\mathbf{y} \mid \mathbf{x}) := p_{T,c}(\mathbf{y} \mid \mathbf{x})$, the transition density of the d -dimensional (double) Langevin (DL) diffusion processes $\mathbf{X}_t^{(c)}$ for $c = 1, \dots, C$, from \mathbf{x} to \mathbf{y} for a pre-defined time $T > 0$ given by

$$d\mathbf{X}_t^{(c)} = \Lambda_c \nabla \log f_c \left(\mathbf{x}_t^{(c)} \right) dt + \Lambda_c^{1/2} d\mathbf{W}_t^{(c)}$$

- $\mathbf{W}_t^{(c)}$ is d -dimensional Brownian motion
- Λ_c is the pre-conditioning matrix associated with sub-posterior f_c
- ∇ is the gradient operator over \mathbf{x}
- Has stationary distribution $f_c^2(\mathbf{x})$

Rejection Sampling (Double Langevin Approach)

- There are many possible choices for $p_c(\mathbf{y} \mid \mathbf{x})$
- Let $p_c(\mathbf{y} \mid \mathbf{x}) := p_{T,c}(\mathbf{y} \mid \mathbf{x})$, the transition density of the d -dimensional (double) Langevin (DL) diffusion processes $\mathbf{X}_t^{(c)}$ for $c = 1, \dots, C$, from \mathbf{x} to \mathbf{y} for a pre-defined time $T > 0$ given by

$$d\mathbf{X}_t^{(c)} = \Lambda_c \nabla \log f_c \left(\mathbf{x}_t^{(c)} \right) dt + \Lambda_c^{1/2} d\mathbf{W}_t^{(c)}$$

- $\mathbf{W}_t^{(c)}$ is d -dimensional Brownian motion
- Λ_c is the pre-conditioning matrix associated with sub-posterior f_c
- ∇ is the gradient operator over \mathbf{x}
- Has stationary distribution $f_c^2(\mathbf{x})$

Rejection Sampling (Double Langevin Approach)

- Extended Target Density:

$$g(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C \left[f_c^2(\mathbf{x}^{(c)}) \cdot p_c(\mathbf{y} | \mathbf{x}^{(c)}) \cdot \frac{1}{f_c(\mathbf{y})} \right]$$

- Consider the proposal density h for the extended target g :

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp \left\{ -\frac{1}{2} (\mathbf{y} - \tilde{\mathbf{x}})^\top \Lambda^{-1} (\mathbf{y} - \tilde{\mathbf{x}}) \right\}$$

- $\tilde{\mathbf{x}} := \left(\sum_{c=1}^C \Lambda_c^{-1} \right)^{-1} \left(\sum_{c=1}^C \Lambda_c^{-1} \mathbf{x}^{(c)} \right)$
- $\Lambda^{-1} := \frac{1}{T} \sum_{c=1}^C \Lambda_c^{-1}$
- T is an arbitrary positive constant

Rejection Sampling (Double Langevin Approach)

- Extended Target Density:

$$g(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C \left[f_c^2(\mathbf{x}^{(c)}) \cdot p_c(\mathbf{y} | \mathbf{x}^{(c)}) \cdot \frac{1}{f_c(\mathbf{y})} \right]$$

- Consider the proposal density h for the extended target g :

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp \left\{ -\frac{1}{2} (\mathbf{y} - \tilde{\mathbf{x}})^\top \Lambda^{-1} (\mathbf{y} - \tilde{\mathbf{x}}) \right\}$$

- $\tilde{\mathbf{x}} := \left(\sum_{c=1}^C \Lambda_c^{-1} \right)^{-1} \left(\sum_{c=1}^C \Lambda_c^{-1} \mathbf{x}^{(c)} \right)$
- $\Lambda^{-1} := \frac{1}{T} \sum_{c=1}^C \Lambda_c^{-1}$
- T is an arbitrary positive constant

Rejection Sampling (Double Langevin Approach)

- Simulation from h is easy:

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp \left\{ -\frac{1}{2}(\mathbf{y} - \tilde{\mathbf{x}})^\top \Lambda^{-1}(\mathbf{y} - \tilde{\mathbf{x}}) \right\}$$

1. Simulate $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$ independently
2. Simulate $\mathbf{y} \sim \mathcal{N}(\tilde{\mathbf{x}}, \Lambda)$

Rejection Sampling (Double Langevin Approach)

- Simulation from h is easy:

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp \left\{ -\frac{1}{2}(\mathbf{y} - \tilde{\mathbf{x}})^\top \Lambda^{-1}(\mathbf{y} - \tilde{\mathbf{x}}) \right\}$$

1. Simulate $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$ independently
2. Simulate $\mathbf{y} \sim \mathcal{N}(\tilde{\mathbf{x}}, \Lambda)$

- └ The Monte Carlo Fusion algorithm
 - └ Constructing a rejection sampler

Rejection Sampling (Double Langevin Approach)

- Simulation from h is easy:

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp \left\{ -\frac{1}{2}(\mathbf{y} - \tilde{\mathbf{x}})^\top \Lambda^{-1}(\mathbf{y} - \tilde{\mathbf{x}}) \right\}$$

1. Simulate $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$ independently
2. Simulate $\mathbf{y} \sim \mathcal{N}(\tilde{\mathbf{x}}, \Lambda)$

Rejection Sampling - acceptance probability

- Acceptance probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho \cdot Q$$

where

$$\left\{ \begin{array}{l} \rho(\mathbf{x}^{(1:C)}) = \exp \left\{ - \sum_{c=1}^C \frac{(\tilde{\mathbf{x}} - \mathbf{x}^{(c)})^\top \Lambda_c^{-1} (\tilde{\mathbf{x}} - \mathbf{x}^{(c)})}{2T} \right\} \\ Q(\mathbf{x}^{(1:C)}, \mathbf{y}) := \prod_{c=1}^C \mathbb{E}_{\mathbb{W}_{\Lambda_c}} \left[\exp \left\{ - \int_0^T \left(\phi_c \left(\mathbf{x}_t^{(c)} \right) - \Phi_c \right) dt \right\} \right] \end{array} \right.$$

where \mathbb{W}_{Λ_c} denotes the law of a Brownian bridge $\{\mathbf{x}_t^{(c)}, t \in [0, T]\}$ with $\mathbf{x}_0^{(c)} := \mathbf{x}^{(c)}$ and $\mathbf{x}_T^{(c)} := \mathbf{y}$ and covariance matrix Λ_c

Rejection Sampling - acceptance probability

- Acceptance probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho \cdot Q$$

where

$$\left\{ \begin{array}{l} \rho(\mathbf{x}^{(1:C)}) = \exp \left\{ - \sum_{c=1}^C \frac{(\tilde{\mathbf{x}} - \mathbf{x}^{(c)})^\top \Lambda_c^{-1} (\tilde{\mathbf{x}} - \mathbf{x}^{(c)})}{2T} \right\} \\ Q(\mathbf{x}^{(1:C)}, \mathbf{y}) := \prod_{c=1}^C \mathbb{E}_{\mathbb{W}_{\Lambda_c}} \left[\exp \left\{ - \int_0^T \left(\phi_c \left(\mathbf{x}_t^{(c)} \right) - \Phi_c \right) dt \right\} \right] \end{array} \right.$$

where \mathbb{W}_{Λ_c} denotes the law of a Brownian bridge $\{\mathbf{x}_t^{(c)}, t \in [0, T]\}$ with $\mathbf{x}_0^{(c)} := \mathbf{x}^{(c)}$ and $\mathbf{x}_T^{(c)} := \mathbf{y}$ and covariance matrix Λ_c

Q Acceptance Probability

$$Q := \prod_{c=1}^C \mathbb{E}_{\mathbb{W}, \Lambda_c} \left[\exp \left\{ - \int_0^T \left(\phi_c \left(\mathbf{x}_t^{(c)} \right) - \Phi_c \right) dt \right\} \right]$$

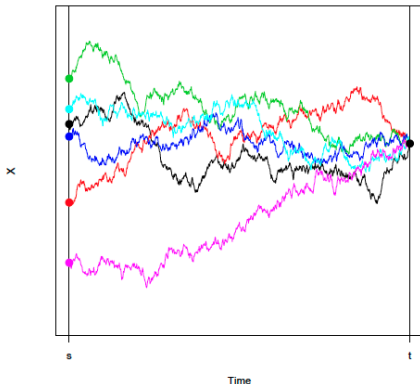
where

$$\phi_c(\mathbf{x}) = \frac{1}{2} \left(\nabla \log f_c(\mathbf{x})^\top \Lambda_c \nabla \log f_c(\mathbf{x}) + \sum_{k=1}^d \Lambda_{c,kk} \frac{\partial \nabla \log f_c(\mathbf{x})}{\partial x_k} \right)$$

- Φ_c are constants such that for all \mathbf{x} , $\phi_c(\mathbf{x}) \geq \Phi_c$ for $c \in \{1, \dots, C\}$
- Events of probability Q can be simulated using **Poisson thinning** and methodology called **Path-space Rejection Sampling (PSRS)** or the **Exact Algorithm** (Beskos et al. [2005], Beskos et al. [2006], Pollock et al. [2016])

Interpretation

- Correct a simple weighted average \tilde{x} of sub-posterior values to a Monte Carlo draw from $\pi(x)$ with acceptance probability $\rho \cdot Q$



Double Langevin Approach - Summary

- Proposal:

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left\{-\frac{1}{2}(\mathbf{y} - \tilde{\mathbf{x}})^T \Lambda^{-1}(\mathbf{y} - \tilde{\mathbf{x}})\right\}$$

- Accept \mathbf{y} as a draw from fusion density π with probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho \cdot Q$$

- Monte Carlo Fusion Algorithm:
 1. Simulate $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$ and $\mathbf{y} \sim \mathcal{N}(\tilde{\mathbf{x}}, \Lambda)$
 2. Accept \mathbf{y} with probability $\rho \cdot Q$

Double Langevin Approach - Summary

- Proposal:

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left\{-\frac{1}{2}(\mathbf{y} - \tilde{\mathbf{x}})^T \Lambda^{-1}(\mathbf{y} - \tilde{\mathbf{x}})\right\}$$

- Accept \mathbf{y} as a draw from fusion density π with probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho \cdot Q$$

- Monte Carlo Fusion Algorithm:
 1. Simulate $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$ and $\mathbf{y} \sim \mathcal{N}(\tilde{\mathbf{x}}, \Lambda)$
 2. Accept \mathbf{y} with probability $\rho \cdot Q$

Double Langevin Approach - Summary

- Proposal:

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left\{-\frac{1}{2}(\mathbf{y} - \tilde{\mathbf{x}})^T \Lambda^{-1}(\mathbf{y} - \tilde{\mathbf{x}})\right\}$$

- Accept \mathbf{y} as a draw from fusion density π with probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho \cdot Q$$

- Monte Carlo Fusion Algorithm:
 1. Simulate $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$ and $\mathbf{y} \sim \mathcal{N}(\tilde{\mathbf{x}}, \Lambda)$
 2. Accept \mathbf{y} with probability $\rho \cdot Q$

- └ The Monte Carlo Fusion algorithm
 - └ Constructing a rejection sampler

Double Langevin Approach - Summary

- Proposal:

$$h(\mathbf{x}^{(1:C)}, \mathbf{y}) \propto \prod_{c=1}^C [f_c(\mathbf{x}^{(c)})] \cdot \exp\left\{-\frac{1}{2}(\mathbf{y} - \tilde{\mathbf{x}})^T \Lambda^{-1}(\mathbf{y} - \tilde{\mathbf{x}})\right\}$$

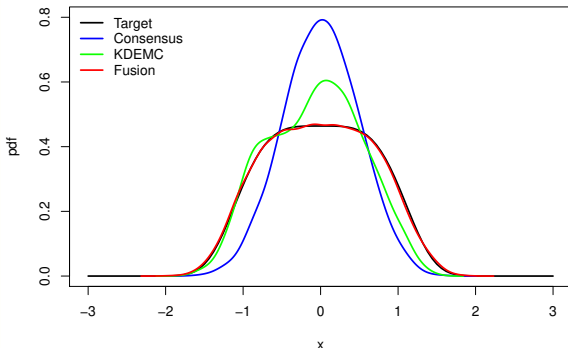
- Accept \mathbf{y} as a draw from fusion density π with probability:

$$\frac{g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})}{h(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}, \mathbf{y})} \propto \rho \cdot Q$$

- Monte Carlo Fusion Algorithm:
 1. Simulate $\mathbf{x}^{(c)} \sim f_c(\mathbf{x})$ and $\mathbf{y} \sim \mathcal{N}(\tilde{\mathbf{x}}, \Lambda)$
 2. Accept \mathbf{y} with probability $\rho \cdot Q$

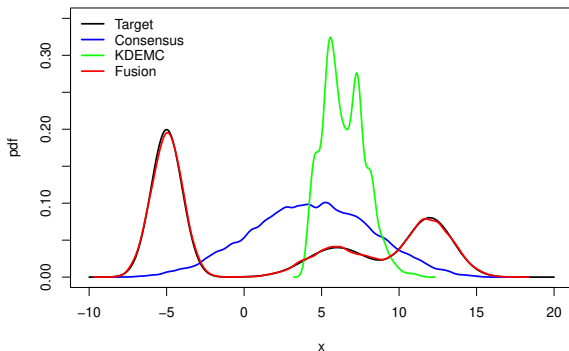
Density with light tails

- **Target:** $\pi(x) \propto e^{-\frac{x^4}{2}}$
- **Sub-posteriors:** $f_c(x) \propto e^{-\frac{x^4}{8}}$ for $c = 1, \dots, 4$
- $N = 20,000$



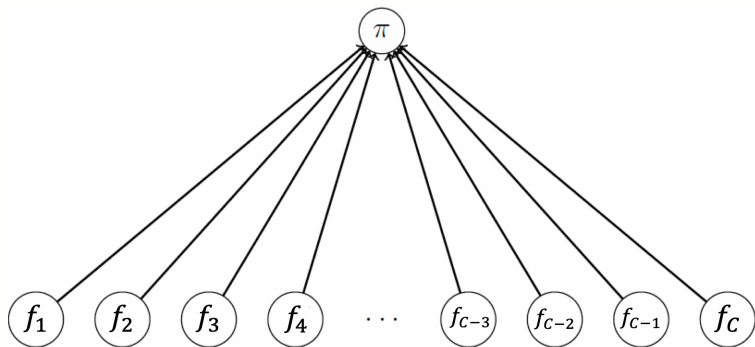
Mixture Gaussian

- **Target:** $\pi(x) \propto 0.5\mathcal{N}(-5, 1) + 0.2\mathcal{N}(6, 2) + 0.3\mathcal{N}(12, 1.5)$
- **Sub-posteriors:** $f_c(x) \propto \pi(x)^{1/4}$ for $c = 1, \dots, 4$
- $N = 20,000$



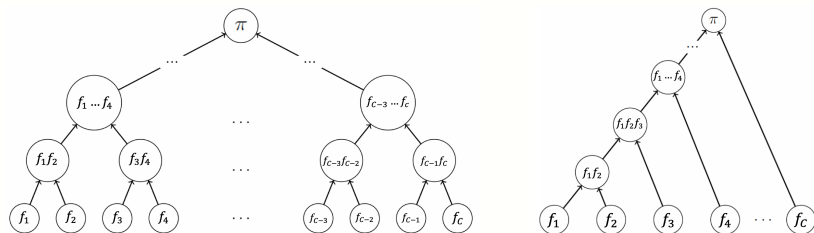
Recall: Fork-and-join

The **fork-and-join** approach:



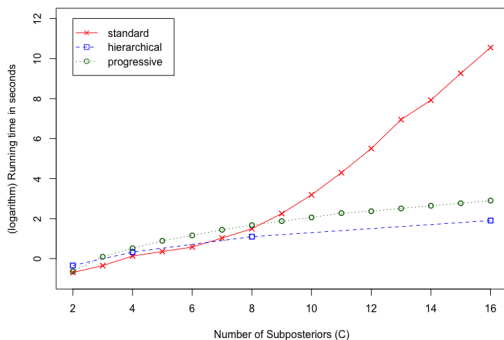
Hierarchical Monte Carlo Fusion

Solution: adopt a **divide-and-conquer** approach:



Example

- Target: $\pi(x) \propto e^{-\frac{x^4}{2}}$
- Sub-posteriors: $f_c(x) = e^{-\frac{x^4}{2c}}$ for $c = 1, \dots, C$



Divide-and-Conquer SMC with Fusion

- Can apply Sequential Monte Carlo in the hierarchical fusion framework
- Rejection sampling can be **wasteful**: large number of proposed samples are rejected
- Motives the use of **Sequential Importance Sampling / Resampling** ideas
 - **Replace** the rejection sampling steps with importance sampling steps
 - Introduce resampling at the nodes if the ESS falls below some threshold
 - Fits into the **Divide-and-Conquer SMC (D&C-SMC)** framework by Lindsten et al. [2017]

Divide-and-Conquer SMC with Fusion

- Can apply Sequential Monte Carlo in the hierarchical fusion framework
- Rejection sampling can be **wasteful**: large number of proposed samples are rejected
- Motives the use of **Sequential Importance Sampling / Resampling** ideas
 - **Replace** the rejection sampling steps with importance sampling steps
 - Introduce resampling at the nodes if the ESS falls below some threshold
 - Fits into the **Divide-and-Conquer SMC (D&C-SMC)** framework by Lindsten et al. [2017]

Divide-and-Conquer SMC with Fusion

- Can apply Sequential Monte Carlo in the hierarchical fusion framework
- Rejection sampling can be **wasteful**: large number of proposed samples are rejected
- Motives the use of **Sequential Importance Sampling / Resampling** ideas
 - **Replace** the rejection sampling steps with importance sampling steps
 - Introduce resampling at the nodes if the ESS falls below some threshold
 - Fits into the **Divide-and-Conquer SMC (D&C-SMC)** framework by Lindsten et al. [2017]

Divide-and-Conquer SMC with Fusion

- Can apply Sequential Monte Carlo in the hierarchical fusion framework
- Rejection sampling can be **wasteful**: large number of proposed samples are rejected
- Motives the use of **Sequential Importance Sampling / Resampling** ideas
 - **Replace** the rejection sampling steps with importance sampling steps
 - Introduce resampling at the nodes if the ESS falls below some threshold
 - Fits into the **Divide-and-Conquer SMC (D&C-SMC)** framework by Lindsten et al. [2017]

Divide-and-Conquer SMC with Fusion

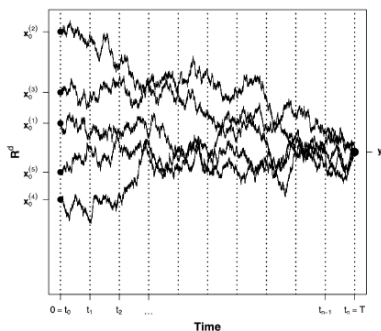
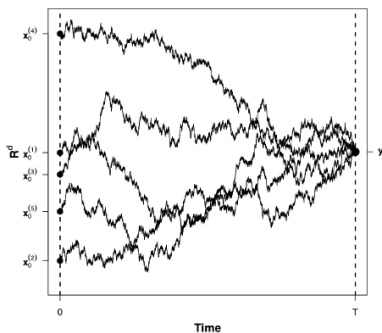
- Can apply Sequential Monte Carlo in the hierarchical fusion framework
- Rejection sampling can be **wasteful**: large number of proposed samples are rejected
- Motives the use of **Sequential Importance Sampling / Resampling** ideas
 - **Replace** the rejection sampling steps with importance sampling steps
 - Introduce resampling at the nodes if the ESS falls below some threshold
 - Fits into the **Divide-and-Conquer SMC (D&C-SMC)** framework by Lindsten et al. [2017]

Divide-and-Conquer SMC with Fusion

- Can apply Sequential Monte Carlo in the hierarchical fusion framework
- Rejection sampling can be **wasteful**: large number of proposed samples are rejected
- Motives the use of **Sequential Importance Sampling / Resampling** ideas
 - **Replace** the rejection sampling steps with importance sampling steps
 - Introduce resampling at the nodes if the ESS falls below some threshold
 - Fits into the **Divide-and-Conquer SMC (D&C-SMC)** framework by Lindsten et al. [2017]

Bayesian Fusion

- Ongoing work by Dai, H., Pollock, M. and Roberts, G.O.
- Tailored to big data Bayesian problems



References

- Beskos, A., Papaspiliopoulos, O., Roberts, G. O., and Fearnhead, P. (2006). Exact and computationally efficient likelihood-based estimation for discretely observed diffusion processes (with discussion). *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(3):333–382.
- Beskos, A., Roberts, G. O., et al. (2005). Exact simulation of diffusions. *The Annals of Applied Probability*, 15(4):2422–2444.
- Dai, H., Pollock, M., and Roberts, G. (2019). Monte Carlo Fusion. *Journal of Applied Probability*, 56(1):174–191.
- Lindsten, F., Johansen, A. M., Naesseth, C. A., Kirkpatrick, B., Schön, T. B., Aston, J., and Bouchard-Côté, A. (2017). Divide-and-conquer with Sequential Monte Carlo. *Journal of Computational and Graphical Statistics*, 26(2):445–458.
- Neiswanger, W., Wang, C., and Xing, E. (2013). Asymptotically exact, embarrassingly parallel MCMC. *arXiv preprint arXiv:1311.4780*.
- Pollock, M., Johansen, A. M., Roberts, G. O., et al. (2016). On the exact and ϵ -strong simulation of (jump) diffusions. *Bernoulli*, 22(2):794–856.
- Scott, S. L., Blocker, A. W., Bonassi, F. V., Chipman, H. A., George, E. I., and McCulloch, R. E. (2016). Bayes and big data: The consensus Monte Carlo algorithm. *International Journal of Management Science and Engineering Management*, 11(2):78–88.