

# MCMC: Intermediate Examples

version: 2020-03-12 · 11:16:30

# Example: Noncentral t-Distribution

**Definition:** Let  $z \sim \mathcal{N}(\mu, \sigma^2)$   $\Pi$   $x \sim \chi^2_{(\nu)}$ . Then

$$y = \frac{z}{\sqrt{x/\nu}} + \eta$$

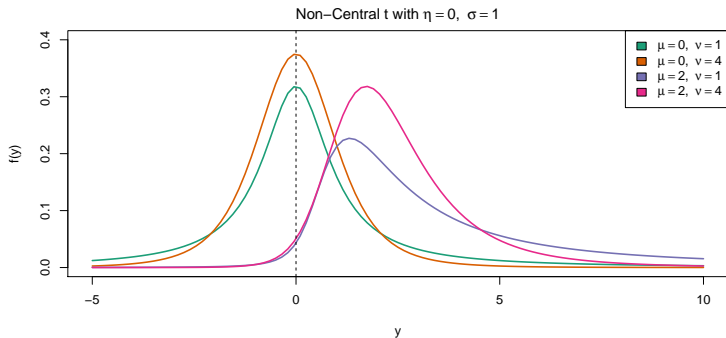
has a Noncentral Student- $t$  distribution, denoted  $y \sim t_{(\nu)}(\mu, \sigma, \eta)$ .

# Noncentral t-Distribution

**Definition:** Let  $z \sim \mathcal{N}(\mu, \sigma^2)$   $\Pi$   $x \sim \chi^2_{(\nu)}$ . Then

$$y = \frac{z}{\sqrt{x/\nu}} + \eta \sim t_{(\nu)}(\mu, \sigma, \eta).$$

**Modeling:** Allows very general specification of mean, variance, skewness and kurtosis.



# Parameter Inference

► **Model:**  $y_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\mu, \sigma, \eta)$

► **Loglikelihood:**

$$\ell(\mu, \sigma, \eta, \nu \mid \mathbf{y}) = \sum_{i=1}^n \text{dt}(x = y_i - \eta) / \sigma, \text{ df} = \nu, \text{ ncp} = \mu, \text{ log} = \text{TRUE}) - n \log \sigma.$$

So for this problem we could get away with MLE, or

# Approximate Bayesian Inference

## 1. Unconstrain Parameters:

$$\theta = (\mu, \sigma, \eta, \nu) \rightarrow \psi = (\mu, \lambda = \log \sigma, \eta, \omega = \log \nu).$$

(Approximation works much better on unconstrained scale.)

# Approximate Bayesian Inference

## 1. Unconstrain Parameters:

$$\boldsymbol{\theta} = (\mu, \sigma, \eta, \nu) \quad \rightarrow \quad \boldsymbol{\psi} = (\mu, \lambda = \log \sigma, \eta, \omega = \log \nu).$$

(Approximation works much better on unconstrained scale.)

## 2. Posterior: $p(\boldsymbol{\psi} \mid \mathbf{y}) \propto \mathcal{L}(\boldsymbol{\psi} \mid \mathbf{y}) \cdot \pi(\boldsymbol{\psi})$ .

# Approximate Bayesian Inference

## 1. Unconstrain Parameters:

$$\boldsymbol{\theta} = (\mu, \sigma, \eta, \nu) \rightarrow \boldsymbol{\psi} = (\mu, \lambda = \log \sigma, \eta, \omega = \log \nu).$$

(Approximation works much better on unconstrained scale.)

## 2. Posterior:

$$p(\boldsymbol{\psi} | \mathbf{y}) \propto \mathcal{L}(\boldsymbol{\psi} | \mathbf{y}) \cdot \pi(\boldsymbol{\psi}).$$

## 3. Normal Approximation:

$$\boldsymbol{\psi} | \mathbf{y} \approx \mathcal{N}(\hat{\boldsymbol{\psi}}, \hat{\mathbf{V}}), \text{ where}$$

$$\hat{\boldsymbol{\psi}} = \arg \max_{\boldsymbol{\psi}} \log p(\boldsymbol{\psi} | \mathbf{y}), \quad \hat{\mathbf{V}} = - \left[ \frac{\partial^2}{\partial \boldsymbol{\psi}^2} \log p(\hat{\boldsymbol{\psi}} | \mathbf{y}) \right]^{-1}.$$

(Also called the mode-quadrature approximation.)

# Approximate Bayesian Inference

## 1. Unconstrain Parameters:

$$\theta = (\mu, \sigma, \eta, \nu) \rightarrow \psi = (\mu, \lambda = \log \sigma, \eta, \omega = \log \nu).$$

(Approximation works much better on unconstrained scale.)

## 2. Posterior: $p(\psi | \mathbf{y}) \propto \mathcal{L}(\psi | \mathbf{y}) \cdot \pi(\psi).$

## 3. Normal Approximation: $\psi | \mathbf{y} \approx \mathcal{N}(\hat{\psi}, \hat{\mathbf{V}})$ , where

$$\hat{\psi} = \arg \max_{\psi} \log p(\psi | \mathbf{y}), \quad \hat{\mathbf{V}} = - \left[ \frac{\partial^2}{\partial \psi^2} \log p(\hat{\psi} | \mathbf{y}) \right]^{-1}.$$

(Also called the mode-quadrature approximation.)

## 4. Monte Carlo Sampling:

- i.  $\psi^{(1)}, \dots, \psi^{(M)} \stackrel{\text{iid}}{\sim} \mathcal{N}(\hat{\psi}, \hat{\mathbf{V}}).$
- ii.  $\theta^{(m)} = (\mu^{(m)}, \exp(\lambda^{(m)}), \eta^{(m)}, \exp(\omega^{(m)})).$



# Parameter Inference

► **Model:**  $y_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\mu, \sigma, \eta)$

► **Loglikelihood:**

$$\ell(\mu, \sigma, \eta, \nu \mid \mathbf{y}) = \sum_{i=1}^n \text{dt}(\mathbf{x} = y_i - \eta) / \sigma, \text{ df} = \nu, \text{ ncp} = \mu, \text{ log} = \text{TRUE}) - n \log \sigma.$$

So for this problem we could get away with MLE, or Approximate Bayesian Inference.

► **However:**

- Don't have gradients for noncentral-t in **TMB**.
- What if we had  $y \mid \mathbf{x} \sim t_{(\nu)}(\mu, \sigma, \mathbf{x}'\boldsymbol{\beta})$ ?

# Parameter Inference

► **Model:**

$$y_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\mu, \sigma, \eta) \iff y_i = \frac{z_i}{\sqrt{x_i/\nu}} + \eta, \quad \begin{aligned} z_i &\stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2) \\ x_i &\stackrel{\text{iid}}{\sim} \chi_{(\nu)}^2 \end{aligned}$$

► **Observed Data:**  $\mathbf{y}_{\text{obs}} = \mathbf{y} = (y_1, \dots, y_n)$ .

► **Missing Data:**  $\mathbf{y}_{\text{miss}} = \mathbf{x} = (x_1, \dots, x_n)$ .

► **Complete Data:**  $\mathbf{y}_{\text{comp}} = (\mathbf{y}, \mathbf{x})$ , with

$$\begin{aligned} x_i &\stackrel{\text{iid}}{\sim} \chi_{(\nu)}^2 \\ y_i | x_i &\stackrel{\text{ind}}{\sim} \mathcal{N}(\eta + \gamma/x_i^{1/2}, \tau^2/x_i), \end{aligned}$$

where  $\gamma = \mu\nu^{1/2}$  and  $\tau = \sigma\nu^{1/2}$ .

# Parameter Inference

- ▶ **Model:**  $y_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\mu, \sigma, \eta)$
- ▶ **Observed Data:**  $\mathbf{y}_{\text{obs}} = \mathbf{y} = (y_1, \dots, y_n)$ .
- ▶ **Complete Data:**  $\mathbf{y}_{\text{comp}} = (\mathbf{y}, \mathbf{x})$ , with

$$\begin{aligned}x_i &\stackrel{\text{iid}}{\sim} \chi_{(\nu)}^2 & \gamma &= \mu\nu^{1/2}, \\y_i | x_i &\stackrel{\text{ind}}{\sim} \mathcal{N}(\eta + \gamma/x_i^{1/2}, \tau^2/x_i), & \tau &= \sigma\nu^{1/2}.\end{aligned}$$

- ▶ **Inference:** Let  $\boldsymbol{\theta} = (\eta, \gamma, \tau^2, \nu)$ .
- ▶ **EM Algorithm:** This would require taking expectations of  $x$ ,  $x^{1/2}$ , and  $\log x$  with respect to

$$\begin{aligned}p(x | y, \boldsymbol{\theta}) &\propto \exp \left\{ -\frac{1}{2} \frac{(y - \eta - \gamma x^{-1/2})^2}{\tau^2/x} + \frac{1}{2} \log x + \left(\frac{\nu-2}{2}\right) \log x - \frac{x}{2} \right\} \\&\propto \exp \{ Ax + Bx^{1/2} + C \log x \},\end{aligned}$$

a nonstandard distribution (don't even know its normalizing constant).

# Parameter Inference

- ▶ **Model:**  $y_i \stackrel{\text{iid}}{\sim} t_{(\nu)}(\mu, \sigma, \eta)$
- ▶ **Observed Data:**  $\mathbf{y}_{\text{obs}} = \mathbf{y} = (y_1, \dots, y_n)$ .
- ▶ **Complete Data:**  $\mathbf{y}_{\text{comp}} = (\mathbf{y}, \mathbf{x})$ , with

$$x_i \stackrel{\text{iid}}{\sim} \chi_{(\nu)}^2,$$
$$y_i | x_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\eta + \gamma/x_i^{1/2}, \tau^2/x_i).$$

- ▶ **Inference:** Let  $\theta = (\eta, \gamma, \tau^2, \nu)$ .
- ▶ **EM Algorithm:** Requires expectations wrt  $p(\mathbf{x} | \mathbf{y}, \theta) \propto \exp \{ A\mathbf{x} + B\mathbf{x}^{1/2} + C \log \mathbf{x} \}$ .
- ▶ **Bayesian Data Augmentation:**

1. Implement an MCMC algorithm on the **augmented** posterior distribution

$$p(\mathbf{x}, \theta | \mathbf{y}) \propto p(\mathbf{y}, \mathbf{x} | \theta) \times \pi(\theta).$$

2. If  $(\mathbf{x}^{(1)}, \theta^{(1)}), \dots, (\mathbf{x}^{(M)}, \theta^{(M)})$  is an MCMC sample from  $p(\mathbf{x}, \theta | \mathbf{y})$ ,  
then the stationary distribution of  $\theta^{(1)}, \dots, \theta^{(M)}$  is  $p(\theta | \mathbf{y}) = \int p(\mathbf{x}, \theta | \mathbf{y}) d\mathbf{x}$ .

(Works for exactly the same reason that the histogram of each random variable in any MCMC converges to its marginal distribution.)

# Bayesian Data Augmentation

- ▶ **Complete Data Likelihood:** Don't cancel out anything involving  $\theta$  or  $\mathbf{x}$ :

$$\begin{aligned}\ell(\theta | \mathbf{x}, \mathbf{y}) &= \log p(\mathbf{y}, \mathbf{x} | \theta) \\ &= -\frac{1}{2} \sum_{i=1}^n \left[ \frac{(y_i - \eta - \gamma x_i^{-1/2})^2}{\tau^2/x_i} - (\nu - 1) \log x_i + x_i \right] \\ &\quad - n \left[ \frac{\tau^2 + \nu}{2} + \log \Gamma(\nu/2) \right].\end{aligned}$$

- ▶ **MCMC Algorithm:** A block Metropolis-within-Gibbs sampler with the following conditional updates:
  - ▶ **Update for  $(\eta, \gamma, \tau)$ :** Canceling everything that doesn't depend on  $\beta = (\eta, \gamma)$  and  $\tau$ , conditional likelihood  $\ell(\beta, \tau | \nu, \mathbf{x}, \mathbf{y})$  is that of a [regression-like](#) model

$$y_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{u}_i' \beta, \tau^2/x_i), \quad \mathbf{u}_i = (1, 1/x_i^{1/2}).$$

# Bayesian Data Augmentation

## ► Complete Data Likelihood:

$$\ell(\theta | \mathbf{x}, \mathbf{y}) = -\frac{1}{2} \sum_{i=1}^n \left[ \frac{(y_i - \eta - \gamma x_i^{-1/2})^2}{\tau^2/x_i} - (\nu - 1) \log x_i + x_i \right] - n \left[ \frac{\tau^2 + \nu}{2} + \log \Gamma \left( \frac{\nu}{2} \right) \right].$$

## ► MCMC Algorithm: A block Metropolis-within-Gibbs sampler with:

- **Update for  $(\eta, \gamma, \tau)$ :** Canceling everything that doesn't depend on  $\beta = (\eta, \gamma)$  and  $\tau$ , conditional likelihood  $\ell(\beta, \tau | \nu, \mathbf{x}, \mathbf{y})$  is that of a **regression-like** model

$$y_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{u}_i' \beta, \tau^2/x_i), \quad \mathbf{u}_i = (1, 1/x_i^{1/2}).$$

- **Conjugate Prior:** Multivariate Normal Inverse-Gamma (mNIX) distribution

$$(\beta, \tau^2) \sim \text{mNIX}(\boldsymbol{\lambda}, \boldsymbol{\Sigma}, \alpha, \gamma) \iff \begin{aligned} \tau^2 &\sim \text{Inv-Gamma}(\alpha, \gamma) \\ \beta | \tau^2 &\sim \mathcal{N}(\boldsymbol{\lambda}, \tau^2 \cdot \boldsymbol{\Sigma}). \end{aligned}$$

$\implies$  **Exact** Gibbs update for  $p(\beta, \tau^2 | \nu, \mathbf{x}, \mathbf{y})$ .

# Bayesian Data Augmentation

## ► Complete Data Likelihood:

$$\ell(\theta | \mathbf{x}, \mathbf{y}) = -\frac{1}{2} \sum_{i=1}^n \left[ \frac{(y_i - \eta - \gamma x_i^{-1/2})^2}{\tau^2/x_i} - (\nu - 1) \log x_i + x_i \right] - n \left[ \frac{\tau^2 + \nu}{2} + \log \Gamma \left( \frac{\nu}{2} \right) \right].$$

## ► MCMC Algorithm: A block Metropolis-within-Gibbs sampler with:

### ► Update for $\nu$ : Conditional likelihood is

$$\ell(\nu | \eta, \gamma, \tau, \mathbf{x}, \mathbf{y}) = -n \log \Gamma(\frac{1}{2}\nu) - \frac{1}{2}\nu \times (n \log(2) - \sum_{i=1}^n \log x_i).$$

### ► Proposal Distribution: Conditional likelihood only depends on $x_i \stackrel{\text{iid}}{\sim} \chi_{(\nu)}^2$ which is an Exponential Family $\implies \ell(\nu | \eta, \gamma, \tau, \mathbf{x}, \mathbf{y})$ is convex. Could do Newton-Raphson to obtain a mode-quadrature normal approximation, but easier to use a random walk proposal.

### ► Prior Distribution: Use $\log \nu \sim \mathcal{N}(0, 2^2)$ . Basically uninformative, since $\Pr(.005 < \nu < 170) \approx 99\%$ (recall that $t_{(\nu=1)} \sim \text{Cauchy}$ and $t_{(\nu \geq 30)} \approx \mathcal{N}(0, 1)$ ). Think of this prior as **regularizing** inference (i.e., prevents $\nu$ from floating off to 0 or $\infty$ ).

# Bayesian Data Augmentation

## ► Complete Data Likelihood:

$$\ell(\theta | \mathbf{x}, \mathbf{y}) = -\frac{1}{2} \sum_{i=1}^n \left[ \frac{(y_i - \eta - \gamma x_i^{-1/2})^2}{\tau^2/x_i} - (\nu - 1) \log x_i + x_i \right] - n \left[ \frac{\tau^2 + \nu}{2} + \log \Gamma \left( \frac{\nu}{2} \right) \right].$$

## ► MCMC Algorithm: A block Metropolis-within-Gibbs sampler with:

### ► Update for $x$ : Conditional posterior is

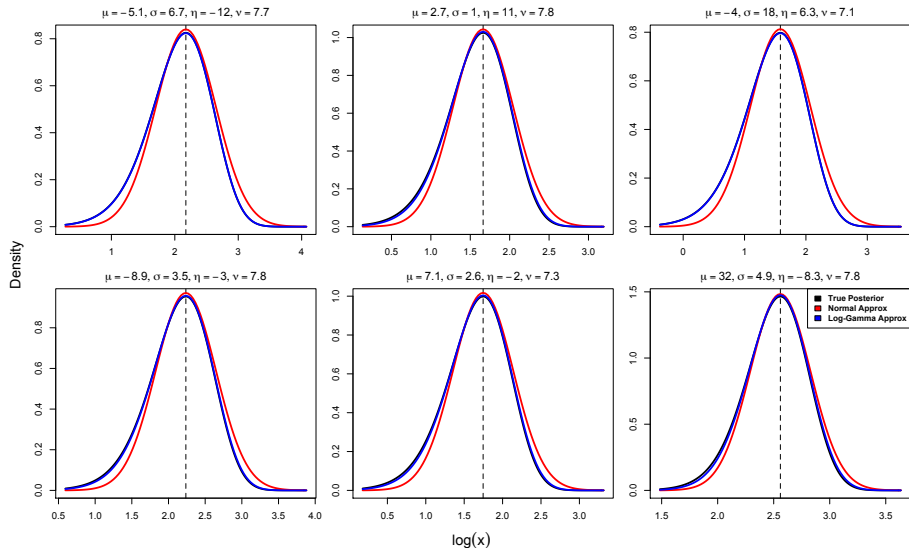
$$p(\mathbf{x} | \mathbf{y}, \theta) \propto \prod_{i=1}^n \exp \left\{ A_i x_i + B_i x_i^{1/2} + C \log x_i \right\}.$$

### ► Proposal Distribution:

- Note that the  $x_i$  are conditionally independent given everything else  $\implies$  exact Gibbs sampler produces IID samples.
- Could do MWG, but this requires  $n$  tuning parameters (one for each  $x_i$ ).
- Note that mode of  $Ax + Bx^{1/2} + C \log x$  has an **analytic** solution  $\implies$  tuning-free MIID-within-Gibbs mode-quadrature proposal.



# Proposal Distribution for $p(x | y, \theta)$



# MCMC Code Checking

- ▶ Much more difficult than checking that  $\hat{\theta} = \arg \max_{\theta} \ell(\theta | \mathbf{y})$ , since
  - ▶ MCMC is a random algorithm
  - ▶ Don't know much about  $p(\theta | \mathbf{y})$  – that's why we're doing MCMC in the first place!
- ▶ **Recommendation:** check code meticulously at every step.

Whenever I skip a step, 99% of time there will be an error and then I don't know if it's in the last step or the one(s) I skipped. So I end up checking every step anyway, except now it takes longer.

# Code Checking Strategies

1. Compare every *simplified* conditional likelihood  $\ell(\theta_j | \boldsymbol{\theta}_{-j}, \mathbf{y})$  to the *unsimplified* likelihood  $\log p(\mathbf{y} | \boldsymbol{\theta})$ .

Difference between the two for any value of  $\theta_j$  should be equal to a constant (possibly depending on  $\mathbf{y}$  and  $\boldsymbol{\theta}_{-j}$ ).

2. Compare every simplified posterior  $p(\theta_j | \boldsymbol{\theta}_{-j}, \mathbf{y})$  to the unsimplified posterior  $\mathcal{L}(\boldsymbol{\theta} | \mathbf{y}) \times \pi(\boldsymbol{\theta})$ .

Same as for loglikelihoods, but now checking Jacobians, i.e., if prior is  $\pi(\boldsymbol{\theta})$  but sampling is done on  $\boldsymbol{\psi} = \mathbf{g}(\boldsymbol{\theta})$ , then  $\pi(\boldsymbol{\psi}) = \pi(\mathbf{g}^{-1}(\boldsymbol{\psi})) \left| \frac{\partial}{\partial \boldsymbol{\psi}} \mathbf{g}^{-1}(\boldsymbol{\psi}) \right|$ .

3. Compare *sampling* from  $p(\theta_j | \boldsymbol{\theta}_{-j}, \mathbf{y})$  to analytic conditional.

To get analytic conditional, recall that  $p(\theta_j | \boldsymbol{\theta}_{-j}, \mathbf{y}) \propto \mathcal{L}(\boldsymbol{\theta} | \mathbf{y}) \times \pi(\boldsymbol{\theta})$ , to normalize evaluate 1-d integral numerically.

4. Compare sampling from  $p(\boldsymbol{\theta} | \mathbf{y})$  for given MCMC to sample from same posterior with a different MCMC.

Both samplers should give same results.

# Code Checking for Noncentral t

**Notation:**  $\theta = (\mu, \sigma, \eta, \nu)$ ,  $\varphi = (\eta, \gamma = \mu\nu^{1/2}, \tau^2 = \sigma^2\nu, \nu) = (\beta, \tau^2, \nu)$ .

## 1. Simplified vs unsimplified likelihoods:

$\ell(\eta, \gamma, \tau^2 | \nu, \mathbf{x}, \mathbf{y})$ ,  $\ell(\nu | \eta, \gamma, \tau^2, \mathbf{x}, \mathbf{y})$ ,  $\log p(\mathbf{x} | \varphi, \mathbf{y})$  can each be checked against

$$p(\mathbf{y}, \mathbf{x} | \varphi) = \underbrace{p(\mathbf{y} | \mathbf{x}, \eta, \gamma, \tau^2)}_{\text{ind} \sim \mathcal{N}(\eta + \gamma \mathbf{x}^{-1/2}, \tau^2 \mathbf{x}^{-1})} \times \underbrace{p(\mathbf{x} | \nu)}_{\text{iid} \sim \chi^2(\nu)}$$

# Code Checking for Noncentral t

**Notation:**  $\theta = (\mu, \sigma, \eta, \nu)$ ,  $\varphi = (\eta, \gamma = \mu\nu^{1/2}, \tau^2 = \sigma^2\nu, \nu) = (\beta, \tau^2, \nu)$ .

## 2. Conditional updates:

►  $p(\nu | \dots)$  and  $p(x_i | \dots)$  compare to analytic 1D posterior  $\propto p(\mathbf{y}, \mathbf{x} | \varphi)\pi(\varphi)$ .

► Prior:  $\log(\nu) \sim \mathcal{N}(\mu_\nu, \sigma_\nu^2)$      $\beta, \tau^2 | \nu \sim \text{mNIX}(\alpha, \gamma, \boldsymbol{\lambda}, \boldsymbol{\Sigma})$

As  $\sigma_\nu, \boldsymbol{\Sigma} \rightarrow \infty$  and  $\alpha, \gamma \rightarrow 0$  this becomes  $\pi(\varphi) \propto 1/\tau^2$

► To check  $p(\beta, \tau^2 | \nu, \mathbf{x}, \mathbf{y}) = \text{mNIX}(\hat{\alpha}, \hat{\gamma}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\Sigma}})$ , note that for any  $\mathbf{a} \in \mathbb{R}^2$ ,

$$\tau^2 | \nu, \mathbf{x}, \mathbf{y} \sim \text{Inv-Gamma}(\hat{\alpha}, \hat{\gamma}), \quad \frac{\mathbf{a}'\beta - \mathbf{a}'\hat{\boldsymbol{\lambda}}}{\sqrt{\hat{\gamma}/\hat{\alpha} \cdot \mathbf{a}'\hat{\boldsymbol{\Sigma}}\mathbf{a}}} | \nu, \mathbf{x}, \mathbf{y} \sim t_{(2\hat{\alpha})}$$

Note that the second result integrates out  $\tau^2$ .

# Code Checking for Noncentral t

**Notation:**  $\theta = (\mu, \sigma, \eta, \nu)$ ,  $\varphi = (\eta, \gamma = \mu\nu^{1/2}, \tau^2 = \sigma^2\nu, \nu) = (\beta, \tau^2, \nu)$ .

## 3. Unconditional Updates:

- ▶ Compare to an MIIID sampler with mode-quadrature normal proposals for  $p(\theta | \mathbf{y}) = p(\mathbf{y} | \theta)\pi(\theta)$ .
- ▶  $p(\mathbf{y} | \theta)$  available through R's built-in function `dt` with `ncp` parameter.
- ▶  $\pi(\theta)$  obtained from  $\pi(\varphi)$  through Jacobian. That is, if  $f_\varphi(\varphi)$  is PDF of prior on  $\varphi$ , then PDF of prior on  $\theta$  is  $f_\theta(\theta) = f_\varphi(\varphi) \times |d\varphi/d\theta|$ , where

$$\frac{d\varphi}{d\theta} = \begin{bmatrix} 0 & \nu^{1/2} & 0 & 0 \\ 0 & 0 & 2\sigma\nu & 0 \\ 1 & 0 & 0 & 0 \\ 0 & \frac{1}{2}\mu\nu^{-1/2} & \sigma^2 & 1 \end{bmatrix} \implies \left| \frac{d\varphi}{d\theta} \right| = 2\sigma\nu^{3/2}.$$

# Code Checking for Noncentral t

**Notation:**  $\theta = (\mu, \sigma, \eta, \nu)$ ,  $\varphi = (\eta, \gamma = \mu\nu^{1/2}, \tau^2 = \sigma^2\nu, \nu) = (\beta, \tau^2, \nu)$ .

## 4. Compare to different MCMC on same posterior:

- ▶ Since this is a 4-parameter problem, probably easiest to compare to MIIID sampling with normal mode-quadrature proposals.
- ▶ For more complicated problems, perhaps easier to use a general-purpose MCMC, which will be slow but easy to program.
- ▶ **Stan:** The state-of-the-art in general-purpose MCMC.
  - ▶ Stan is a programming language very similar to R to which you input an arbitrary  $\log p(\theta | \mathbf{y})$ .
  - ▶ Implements and compiles in C++ a very effective MCMC algorithm called Hybrid Monte Carlo (HMC), but usually referred to as **Hamiltonian Monte Carlo**.

# Hamiltonian Monte Carlo (HMC)

▶ **Problem:** Sample from  $\mathbf{x} \sim p(\mathbf{x}) \propto \exp\{\Omega(\mathbf{x})\}$ ,  $\mathbf{x} = (x_1, \dots, x_d)$ .

▶ **Hamiltonian Dynamics:**

▶ **System Variables:**

▶ *Position Variables*  $\mathbf{x} = (x_1, \dots, x_d)$ .

▶ *Momentum Variables*  $\mathbf{v} = (v_1, \dots, v_d)$ .

▶ *Phase-Space Variables*  $\Gamma = (\mathbf{x}, \mathbf{v})$ .

▶ **Hamiltonian Function:**  $\mathcal{H}(\mathbf{x}, \mathbf{v}) = -\Omega(\mathbf{x}) + \frac{1}{2} \sum_{i=1}^d \frac{v_i^2}{m_i}$ .

▶ **Equations of Motion:** Consider the function  $\Gamma_t = \Gamma(t)$  defined by the system of ordinary differential equations (ODEs) and initial conditions

$$\frac{d}{dt} x_i(t) = \frac{v_i(t)}{m_i}, \quad \frac{d}{dt} v_i(t) = \frac{\partial}{\partial x_i} \Omega(\mathbf{x}_t), \quad \Gamma_0 = (\mathbf{x}_0, \mathbf{v}_0).$$

Thus we have some function  $\Psi : \mathbb{R}^{2d} \times \mathbb{R} \rightarrow \mathbb{R}^{2d}$  such that  $\Psi(\Gamma_0, t) = \Gamma_t$ .



# Hamiltonian Monte Carlo (HMC)

► **Problem:** Sample from  $\mathbf{x} \sim p(\mathbf{x}) \propto \exp\{\Omega(\mathbf{x})\}$ ,  $\mathbf{x} = (x_1, \dots, x_d)$ .

► **Hamiltonian Dynamics:**

► **System Variables:**  $\mathbf{x}$  (position),  $\mathbf{v}$  (momentum),  $\Gamma = (\mathbf{x}, \mathbf{v})$  (phase-space).

► **Hamiltonian Function:** 
$$\mathcal{H}(\mathbf{x}, \mathbf{v}) = -\Omega(\mathbf{x}) + \frac{1}{2} \sum_{i=1}^d \frac{v_i^2}{m_i}.$$

► **Equations of Motion:** Define  $\Gamma_t = \Psi(\Gamma_0, t)$  as solution to system of ODEs

$$\frac{d}{dt} x_i(t) = \frac{v_i(t)}{m_i}, \quad \frac{d}{dt} v_i(t) = \frac{\partial}{\partial x_i} \Omega(\mathbf{x}_t), \quad \Gamma_0 = (\mathbf{x}_0, \mathbf{v}_0).$$

► **Conservation of Energy:** If  $\Gamma_t = \Psi(\Gamma_0, t)$ , then  $\mathcal{H}(\mathbf{x}_0, \mathbf{v}_0) = \mathcal{H}(\mathbf{x}_t, \mathbf{v}_t)$ .

► **Preservation of Volume:** Change of variables  $\Gamma_0 \rightarrow \Gamma_t$  has Jacobian

$$\left| \frac{d\Gamma_t}{d\Gamma_0} \right| = \left| \frac{d}{d\Gamma_0} \Psi(\Gamma_0, t) \right| = 1.$$

# Hamiltonian Monte Carlo (HMC)

► **Problem:** Sample from  $\mathbf{x} \sim p(\mathbf{x}) \propto \exp\{\Omega(\mathbf{x})\}$ ,  $\mathbf{x} = (x_1, \dots, x_d)$ .

► **Hamiltonian Dynamics:**

► **Hamiltonian Function:** 
$$\mathcal{H}(\mathbf{x}, \mathbf{v}) = -\Omega(\mathbf{x}) + \frac{1}{2} \sum_{i=1}^d \frac{v_i^2}{m_i}.$$

► **Equations of Motion:** Define  $\Gamma_t = \Psi(\Gamma_0, t)$  as solution to system of ODEs

$$\frac{d}{dt} x_i(t) = \frac{v_i(t)}{m_i}, \quad \frac{d}{dt} v_i(t) = \frac{\partial}{\partial x_i} \Omega(\mathbf{x}_t), \quad \Gamma_0 = (\mathbf{x}_0, \mathbf{v}_0).$$

► **(Idealized) HMC Proposal:** Given  $\mathbf{x}_{\text{curr}}$  and  $L > 0$ :

1. Let  $\Gamma_0 = (\mathbf{x}_{\text{curr}}, \mathbf{v}_0)$ , where  $v_{0i} \stackrel{\text{ind}}{\sim} \mathcal{N}(0, m_i)$ .

2. Let  $\mathbf{x}_{\text{prop}} = \mathbf{x}_L$ , where  $(\mathbf{x}_L, \mathbf{v}_L) = \Gamma_L = \Psi(\Gamma_0, L)$

⇒ Metropolis-Hastings acceptance rate:

# Hamiltonian Monte Carlo (HMC)

► **Problem:** Sample from  $\mathbf{x} \sim p(\mathbf{x}) \propto \exp\{\Omega(\mathbf{x})\}$ ,  $\mathbf{x} = (x_1, \dots, x_d)$ .

► **Hamiltonian Dynamics:**

► **Hamiltonian Function:**  $\mathcal{H}(\mathbf{x}, \mathbf{v}) = -\Omega(\mathbf{x}) + \frac{1}{2} \sum_{i=1}^d \frac{v_i^2}{m_i}$ .

► **Equations of Motion:** Define  $\Gamma_t = \Psi(\Gamma_0, t)$  as solution to system of ODEs

$$\frac{d}{dt} x_i(t) = \frac{v_i(t)}{m_i}, \quad \frac{d}{dt} v_i(t) = \frac{\partial}{\partial x_i} \Omega(\mathbf{x}_t), \quad \Gamma_0 = (\mathbf{x}_0, \mathbf{v}_0).$$

► **(Idealized) HMC Proposal:** Given  $\mathbf{x}_{\text{curr}}$  and  $L > 0$ :

1. Let  $\Gamma_0 = (\mathbf{x}_{\text{curr}}, \mathbf{v}_0)$ , where  $v_{0i} \stackrel{\text{ind}}{\sim} \mathcal{N}(0, m_i)$ .

2. Let  $\mathbf{x}_{\text{prop}} = \mathbf{x}_L$ , where  $(\mathbf{x}_L, \mathbf{v}_L) = \Gamma_L = \Psi(\Gamma_0, L)$

⇒ Metropolis-Hastings acceptance rate: 100%!!!

# Hamiltonian Monte Carlo (HMC)

- ▶ **Problem:** Sample from  $\mathbf{x} \sim p(\mathbf{x}) \propto \exp\{\Omega(\mathbf{x})\}$ ,  $\mathbf{x} = (x_1, \dots, x_d)$ .
- ▶ **(Idealized) HMC Proposal:** Given  $\mathbf{x}_{\text{curr}}$  and  $L > 0$ :
  1. Let  $\Gamma_0 = (\mathbf{x}_{\text{curr}}, \mathbf{v}_0)$ , where  $v_{0i} \stackrel{\text{ind}}{\sim} \mathcal{N}(0, m_i)$ .
  2. Let  $\mathbf{x}_{\text{prop}} = \mathbf{x}_L$ , where  $(\mathbf{x}_L, \mathbf{v}_L) = \Gamma_L = \Psi(\Gamma_0, L)$   
 $\implies$  Metropolis-Hastings acceptance rate: 100%!!!
- ▶ **In Practice:**
  - ▶ Can't solve ODE exactly: discretize  $\implies$  acceptance rate  $\neq 1$ .
  - ▶ Lots of tuning parameters: ODE solver step size, total time  $t$ , mass  $m$ .
  - ▶ Gradients: To solve ODE need  $\frac{\partial}{\partial x_i} \Omega(\mathbf{x})$ , which is a lot of programming effort.

# Hamiltonian Monte Carlo (HMC)

► **Problem:** Sample from  $\mathbf{x} \sim p(\mathbf{x}) \propto \exp\{\Omega(\mathbf{x})\}$ ,  $\mathbf{x} = (x_1, \dots, x_d)$ .

► **(Idealized) HMC Proposal:** Given  $\mathbf{x}_{\text{curr}}$  and  $L > 0$ :

1. Let  $\Gamma_0 = (\mathbf{x}_{\text{curr}}, \mathbf{v}_0)$ , where  $v_{0i} \stackrel{\text{ind}}{\sim} \mathcal{N}(0, m_i)$ .

2. Let  $\mathbf{x}_{\text{prop}} = \mathbf{x}_L$ , where  $(\mathbf{x}_L, \mathbf{v}_L) = \Gamma_L = \Psi(\Gamma_0, L)$

⇒ Metropolis-Hastings acceptance rate: 100%!!!

► **In Practice:**

► Can't solve ODE exactly: discretize ⇒ acceptance rate  $\neq 1$ .

► Lots of tuning parameters: ODE solver step size, total time  $t$ , mass  $m$ .

► Gradients: To solve ODE need  $\frac{\partial}{\partial x_i} \Omega(\mathbf{x})$ , which is a lot of programming effort.

All of this is done **automatically** by Stan :)

# Stan Examples

## ► Examples:

1. **Curved Mean-Variance Normal:**  $\sigma \sim p(\sigma | \mathbf{y})$ , where

$$\sigma \sim \chi_{(7)}^2, \quad y_i | \sigma \stackrel{\text{iid}}{\sim} \mathcal{N}(\sigma, \sigma^2).$$

2. **Banana distribution:**  $\mathbf{x} \sim p(\mathbf{x} | \sigma, y)$ , where  $\mathbf{x} = (x_1, x_2)$  and

$$p(\mathbf{x} | \sigma, y) \propto \exp \left\{ - \left[ \frac{(y - x_1 \cdot x_2)^2}{2\sigma^2} + \frac{(x_1 - x_2)^2}{2} \right] \right\}$$

## ► Key Concepts:

- **Testing Stan code:** Generic MCMC, so only need to check log-posterior is correct. Do this with **rstan** package functions `log_prob` and `expose_stan_functions`.
- **Testing other code:** Stan is relatively easy to program, so use it to compare to sampling from a more specific MCMC algorithm for a particular problem (can often do better than any generic algorithm at expense of human hours).

# Stan Resources

- ▶ Instructions for installing Stan in R can be found [here](#). Follow these **to the letter** or Stan probably won't work properly!
- ▶ Full Stan documentation (tons of examples) and `rstan` package vignette (for operating Stan from within R) can be found [here](#).
- ▶ Detailed explanation of HMC algorithm, its strengths and pitfalls, and many of its variants can be found [here](#).