

# Bayesian Computation Using Design of Experiments-Based Interpolation Technique

V. Roshan Joseph

*Technometrics*, August, 2012  
(with discussions)

Supported by  
NSF CMMI-1030125

# Bayesian Computation

- Many intractable high-dimensional integrals
  - Posterior distribution
  - Posterior summaries
  - Marginal posterior distributions
  - Posterior predictive distributions
- Use approximation methods

# Deterministic Methods

- Laplace's Approximation (e.g. Tierney and Kadane 1986)
  - May not be accurate.
- Gaussian Quadrature (e.g. Naylor and Smith 1982)
  - Can produce accurate results.
  - Cannot be used for high dimensions.

# Simulation-Based Methods

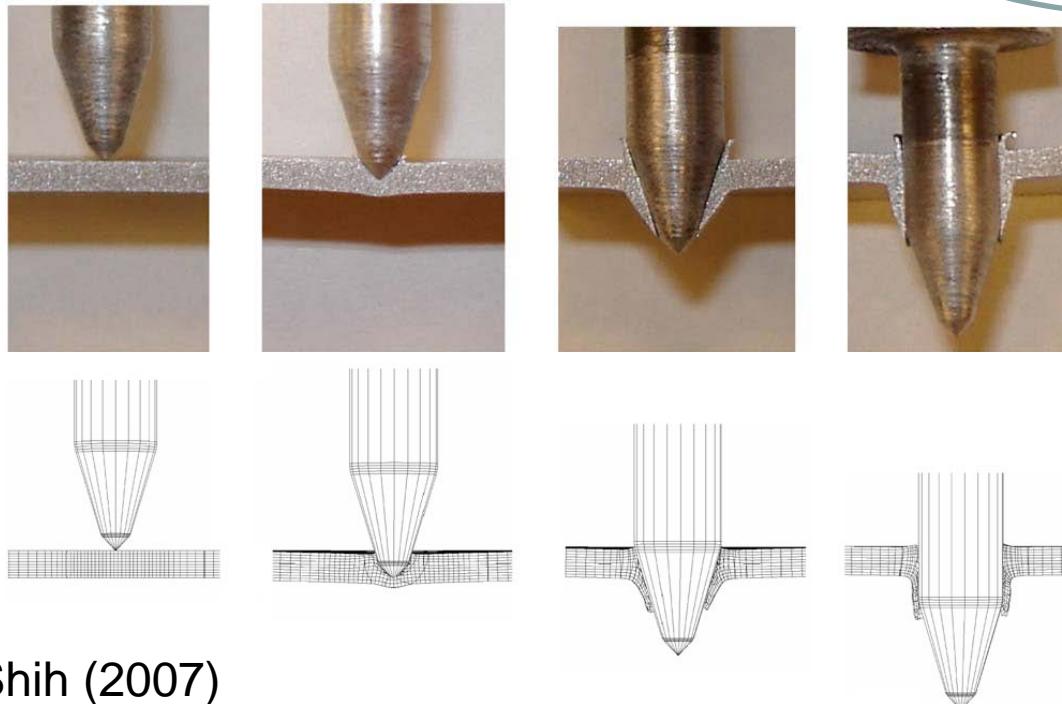
- MC/MCMC (Metropolis et al. 1953, Hastings 1970, Geman and Geman 1984, Gelfand and Smith 1990, ...)
  - Can obtain results with arbitrary precision
  - Suffer less from the curse of dimensionality
  - Convergence issues
  - High computational cost when dealing with computationally expensive posteriors

# Examples of computationally expensive posteriors

- Model calibration

$$\min_{\theta} \sum_{i=1}^n \{y_i - f(x_i; \theta)\}^2$$

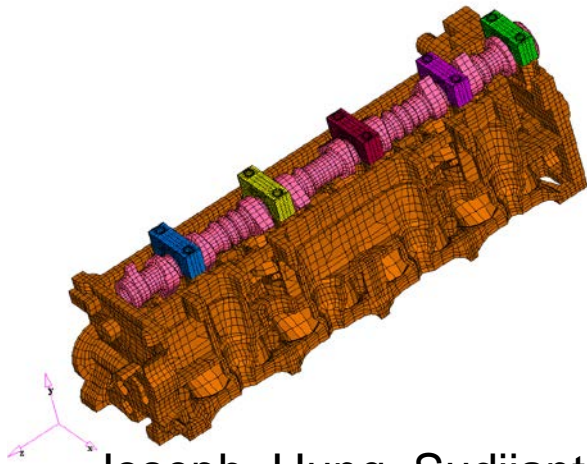
Expensive



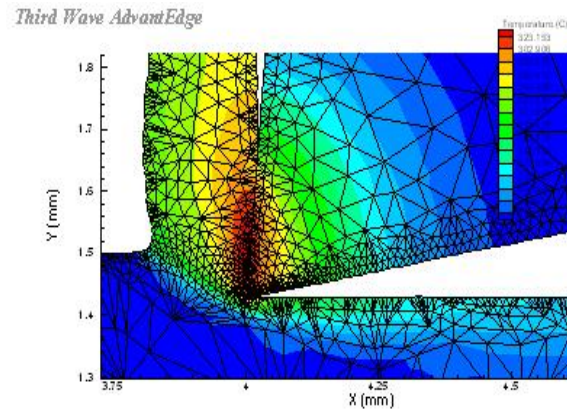
Miller and Shih (2007)

# Examples-continued

- Geostatistics/Spatial statistics/Computer experiments



Joseph, Hung, Sudjianto (2008)



Hung, Joseph, Melkote (2009)

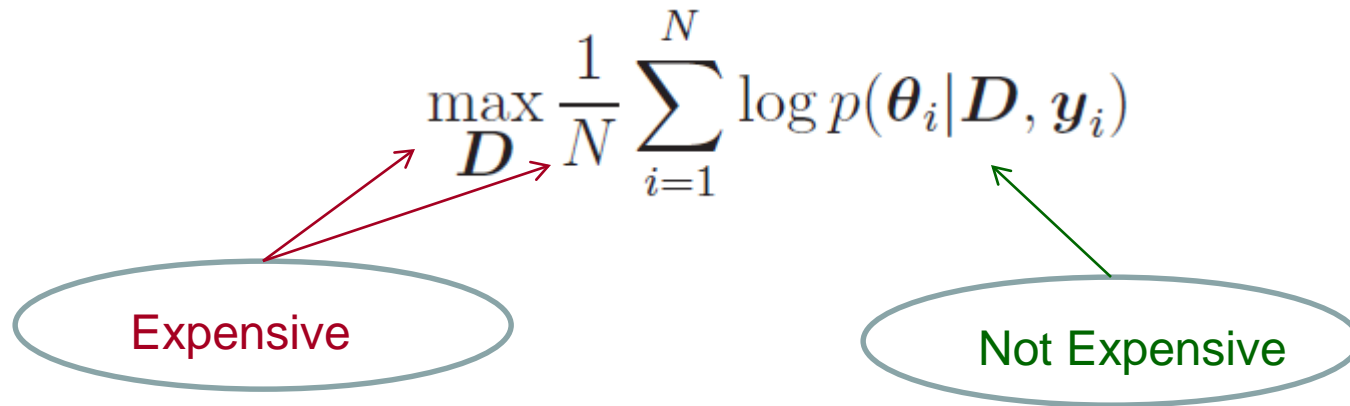
$$f(x)|\mu, \sigma^2, \theta \sim GP(\mu, \sigma^2 R(:, \theta))$$

$$p(y|\mu, \sigma^2, \theta) \propto \frac{1}{\sigma^n |R(\theta)|^{1/2}} \exp\left\{-\frac{1}{2\sigma^2} (y - \mu \mathbf{1})' R(\theta)^{-1} (y - \mu \mathbf{1})\right\}$$

Expensive

# Examples-continued

- Inexpensive likelihood, but appearing within some algorithms.
- Simulation-based nonlinear optimal design (Muller 1999)



# New Deterministic Method

- **Design of Experiments-based Interpolation technique (DoIt)**
  - Design of experiments
  - Interpolation methods (e.g., kriging)
- **Advantages**
  - Can obtain results with arbitrary precision
  - Suffer less from the curse of dimensionality
  - Works much faster than MC/MCMC
- **Disadvantages**
  - Small to moderate number dimensions
  - Continuous parameters

# Earlier Work

- Bayes-Hermite Quadrature
  - O'Hagan (1991)
  - Kennedy (1998)
  - Rasmussen and Ghahramani (2003)
- Hybrid Monte Carlo using Gaussian Process Models
  - Rasmussen (2003)
  - Bliznyuk et al. (2008)
  - Henedrson et al. (2008)
  - Fielding, Nott, and Liong (2011)

# Laplace's Approximation

- Bayesian model:

$$\mathbf{y}|\boldsymbol{\theta} \sim p(\mathbf{y}|\boldsymbol{\theta})$$

$$\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$$

- Unnormalized posterior:  $h(\boldsymbol{\theta}) \propto p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})$

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} h(\boldsymbol{\theta})$$

$$h(\boldsymbol{\theta}) \approx h(\hat{\boldsymbol{\theta}}) \exp\left\{-\frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})' \Sigma^{-1}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})\right\}$$

$$\boldsymbol{\theta}|\mathbf{y} \sim^a N(\hat{\boldsymbol{\theta}}, \Sigma)$$

$$\Sigma = [-\nabla^2 \log(h(\hat{\boldsymbol{\theta}}))]^{-1}$$

# DoIt

- Unnormalized normal density function:

$$g(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \exp\left\{-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu})\right\}$$

- Experimental design:

$$\boldsymbol{D} = \{\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_m\}$$

- DoIt:

$$h(\boldsymbol{\theta}) \approx \sum_{i=1}^m c_i g(\boldsymbol{\theta}; \boldsymbol{\nu}_i, \boldsymbol{\Sigma})$$

- Kriging, Radial Basis Function,...

# DoIt-continued

- Evaluations:  $\mathbf{h} = (h_1, \dots, h_m)'$ , where  $h_i = h(\boldsymbol{\nu}_i)$
- To get interpolation:

$$\mathbf{G}\mathbf{c} = \mathbf{h},$$

where  $\mathbf{c} = (c_1, \dots, c_m)'$  and  $G_{ij} = g(\boldsymbol{\nu}_i; \boldsymbol{\nu}_j, \boldsymbol{\Sigma})$

$$\tilde{\mathbf{c}} = \mathbf{G}^{-1}\mathbf{h}$$

- Let  $\mathbf{g}(\boldsymbol{\theta}) = (g(\boldsymbol{\theta}; \boldsymbol{\nu}_1, \boldsymbol{\Sigma}), \dots, g(\boldsymbol{\theta}; \boldsymbol{\nu}_m, \boldsymbol{\Sigma}))'$

$$\tilde{h}(\boldsymbol{\theta}) = \tilde{\mathbf{c}}'\mathbf{g}(\boldsymbol{\theta})$$

# DoIt-continued

- Marginal likelihood

$$\begin{aligned}\int \tilde{h}(\boldsymbol{\theta}) d\boldsymbol{\theta} &= \tilde{\mathbf{c}}' \int \mathbf{g}(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= (2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2} \tilde{\mathbf{c}}' \mathbf{1}\end{aligned}$$

- Posterior distribution

$$\tilde{p}(\boldsymbol{\theta}|\mathbf{y}) \approx \frac{\tilde{\mathbf{c}}' \mathbf{g}(\boldsymbol{\theta})}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2} \tilde{\mathbf{c}}' \mathbf{1}} = \frac{\tilde{\mathbf{c}}' \boldsymbol{\phi}(\boldsymbol{\theta})}{\tilde{\mathbf{c}}' \mathbf{1}}$$

$$\tilde{p}(\boldsymbol{\theta}|\mathbf{y}) \approx \frac{\sum_{i=1}^m \tilde{c}_i \phi(\boldsymbol{\theta}; \boldsymbol{\nu}_i, \boldsymbol{\Sigma})}{\sum_{i=1}^m \tilde{c}_i}$$

# Example

- Bayesian model:

$$y|\theta \sim \text{Poisson}(\theta),$$

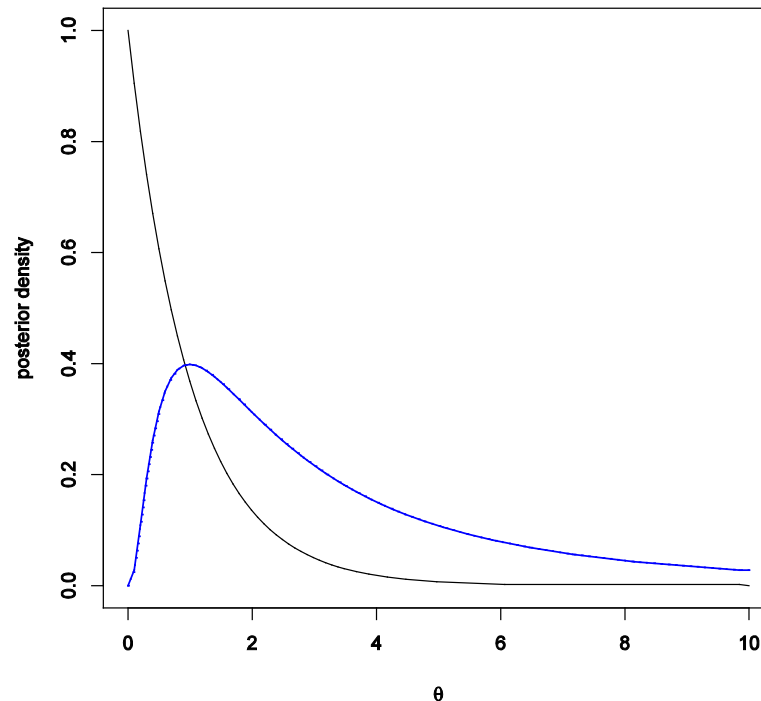
$$p(\theta) \propto 1.$$

- For  $y=0$ , posterior distribution:

$$\theta|y \sim \text{Exp}(1).$$

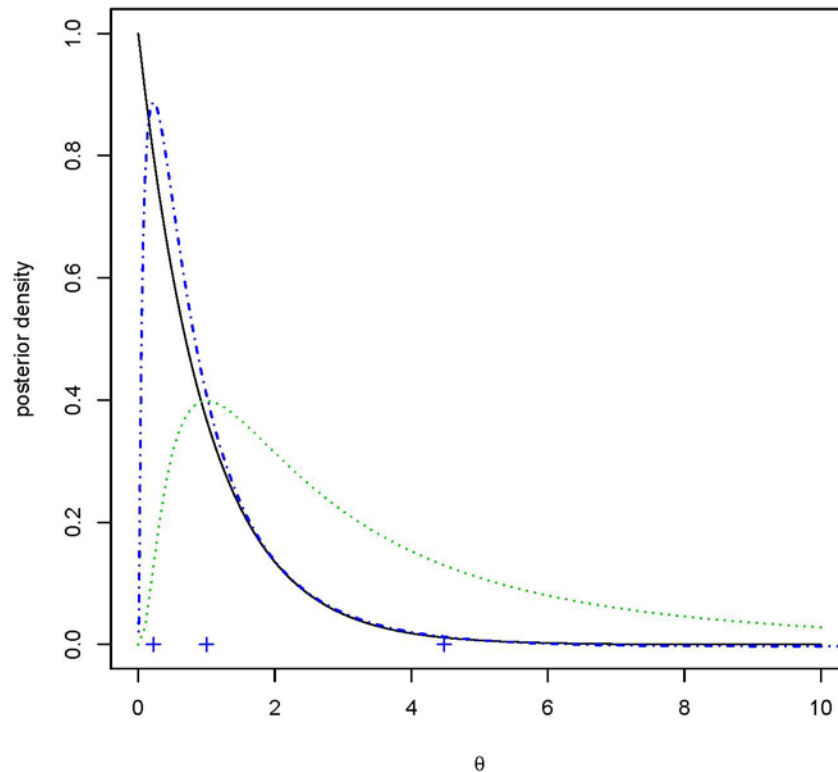
# Example-continued

- Laplace's approximation:  $\gamma = \log(\theta)$
- $\hat{\gamma} = 0$  and  $\hat{\sigma}^2 = 1$ :  $\theta|y \sim^a \log - normal(0, 1)$ .



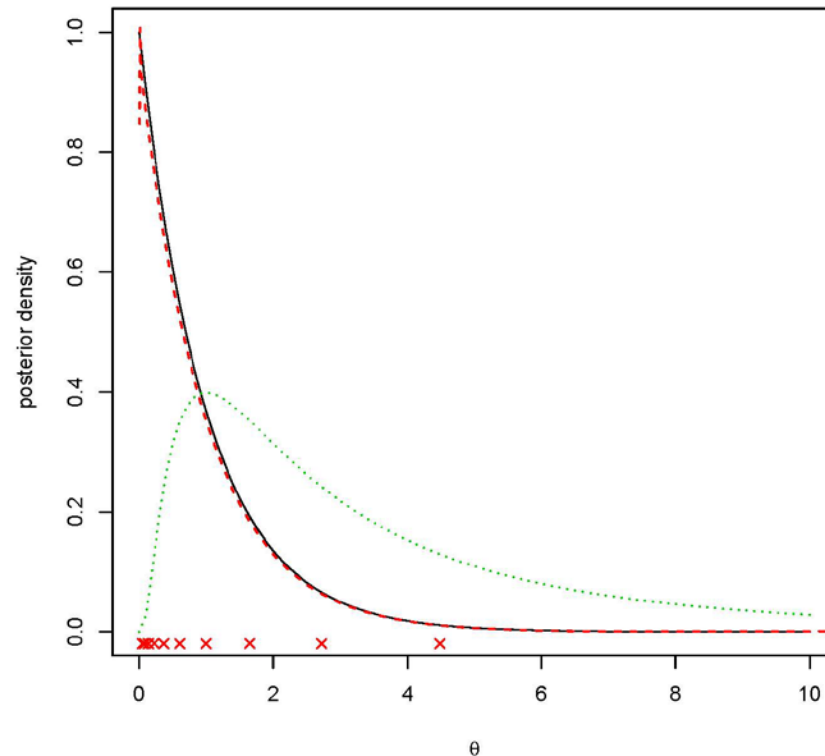
# Example-continued

- DoIt :  $\nu_1 = \hat{\gamma}$ ,  $\nu_2 = \hat{\gamma} - 1.5\sigma$ , and  $\nu_2 = \hat{\gamma} + 1.5\sigma$



# Example-continued

- DoIt : 10 equally spaced points from  $\hat{\gamma} - 3\sigma$  to  $\hat{\gamma} + 1.5\sigma$



# A Result

*Theorem 1:* If  $h(\boldsymbol{\theta})$  is continuous, then for any  $\alpha \in (0, 1)$  and any  $\epsilon > 0$ , there exists a finite number of points  $D = \{\nu_1, \dots, \nu_m\}$  in  $\Theta$  such that

$$\left| \frac{\hat{h}(\boldsymbol{\theta}) / \int_{\Theta} \hat{h}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{h(\boldsymbol{\theta}) / \int_{\Theta} h(\boldsymbol{\theta}) d\boldsymbol{\theta}} - 1 \right| < \epsilon \quad (6)$$

for all  $\boldsymbol{\theta} \in \Theta$ , where  $\hat{h}(\boldsymbol{\theta})$  is any continuous and uniformly convergent interpolator of  $h(\boldsymbol{\theta})$  on  $D$  and  $\Theta$  is the  $(1 - \alpha)$  highest posterior density (HPD) credible set.

As  $\alpha \rightarrow 0$ ,

$$\frac{\hat{h}(\boldsymbol{\theta}) / \int_{\Theta} \hat{h}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{h(\boldsymbol{\theta}) / \int_{\Theta} h(\boldsymbol{\theta}) d\boldsymbol{\theta}} \rightarrow \tilde{p}(\boldsymbol{\theta}|\mathbf{y})/p(\boldsymbol{\theta}|\mathbf{y}),$$

and as  $\epsilon \rightarrow 0$ ,

$$\tilde{p}(\boldsymbol{\theta}|\mathbf{y})/p(\boldsymbol{\theta}|\mathbf{y}) \rightarrow 1.$$

# Unknown Posterior Mode

- Assume that  $\Sigma = \text{diag}\{\sigma_1^2, \dots, \sigma_d^2\}$
- Leave-one-out cross validation:  $e_i = h_i - \tilde{h}_{(i)i}$
- From the kriging literature

$$e_i = \frac{(\mathbf{G}^{-1})_i}{(\mathbf{G}^{-1})_{ii}} \mathbf{h}$$

- Minimize

$$MSCV = \frac{1}{m} \mathbf{e}' \mathbf{e}.$$

- or

$$WMSCV = \frac{1}{m} \mathbf{e}' \text{diag}(\mathbf{G}^{-1}) \mathbf{e}$$

# Example

- Bayesian model:

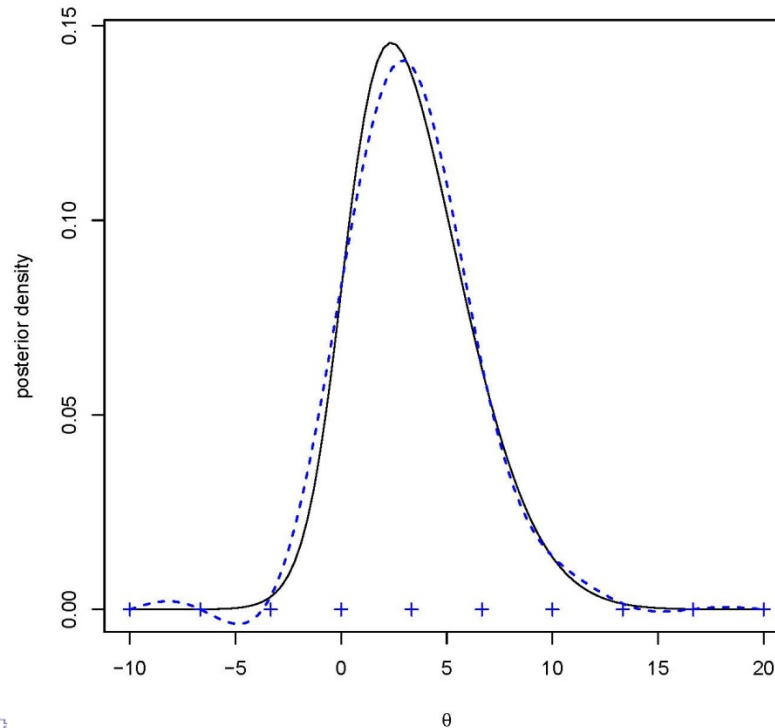
$$y|\theta \sim \text{Bernoulli}(\{1 + \exp(-\theta)\}^{-1}),$$

$$\theta \sim N(\mu, \tau^2).$$

- Suppose  $y = 1$ ,  $\mu = 1$  and  $\tau = 4$ .

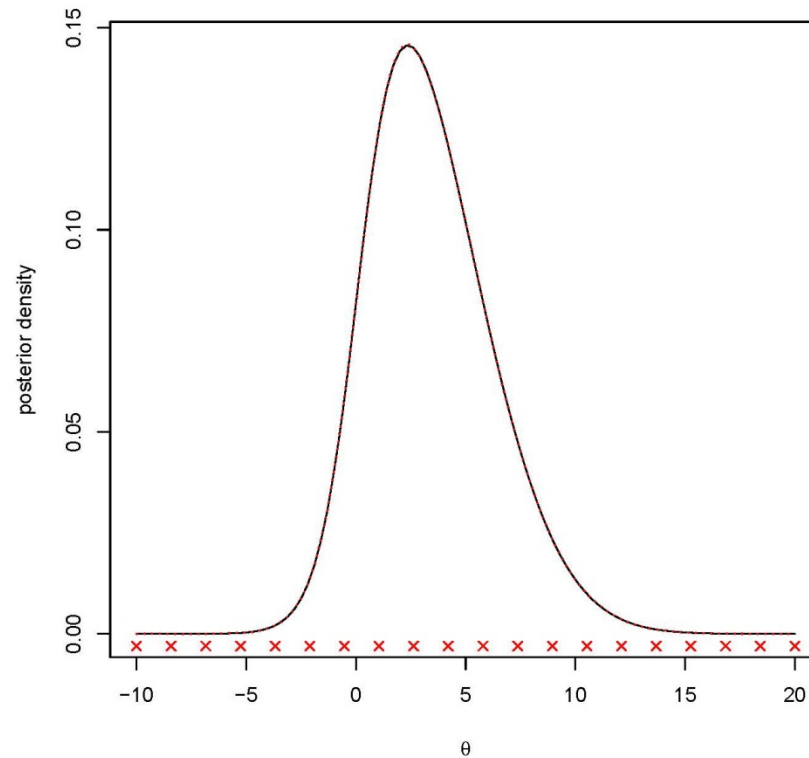
# Example-continued

- Suppose we sample 10 equally spaced points from -10 to 20.
- Minimizing  $WMSCV$ :  $\hat{\sigma}^2 = 9.30$



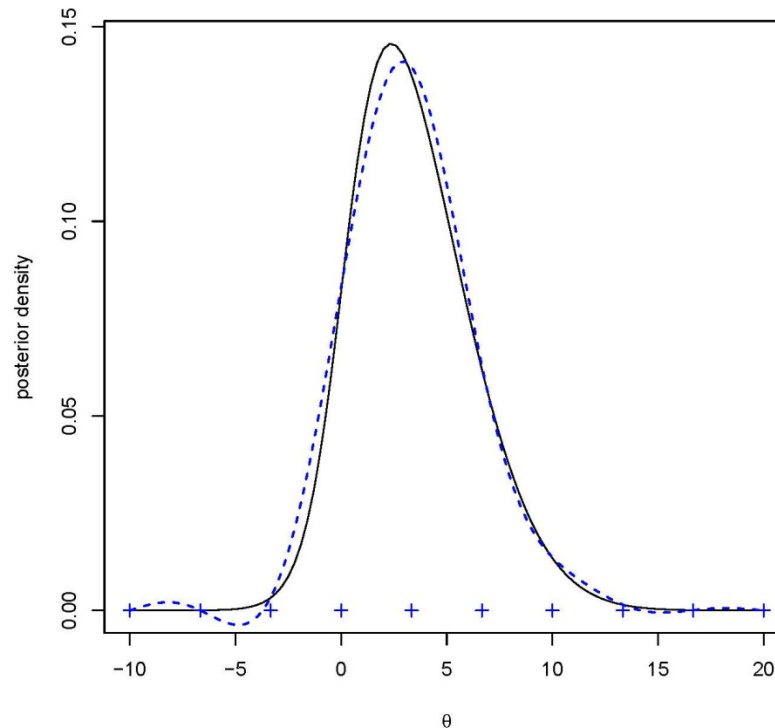
# Example-continued

- $m=20$



# Negativity Problem

- The coefficients  $\tilde{c}_i$  can be negative and can lead to negative posterior density values.



# Mixture Normal Approximation

- Restrict  $c_i$ 's to be nonnegative:

$$\min_{\mathbf{c} \geq \mathbf{0}} (\mathbf{h} - \mathbf{G}\mathbf{c})' \mathbf{G}^{-1} (\mathbf{h} - \mathbf{G}\mathbf{c})$$

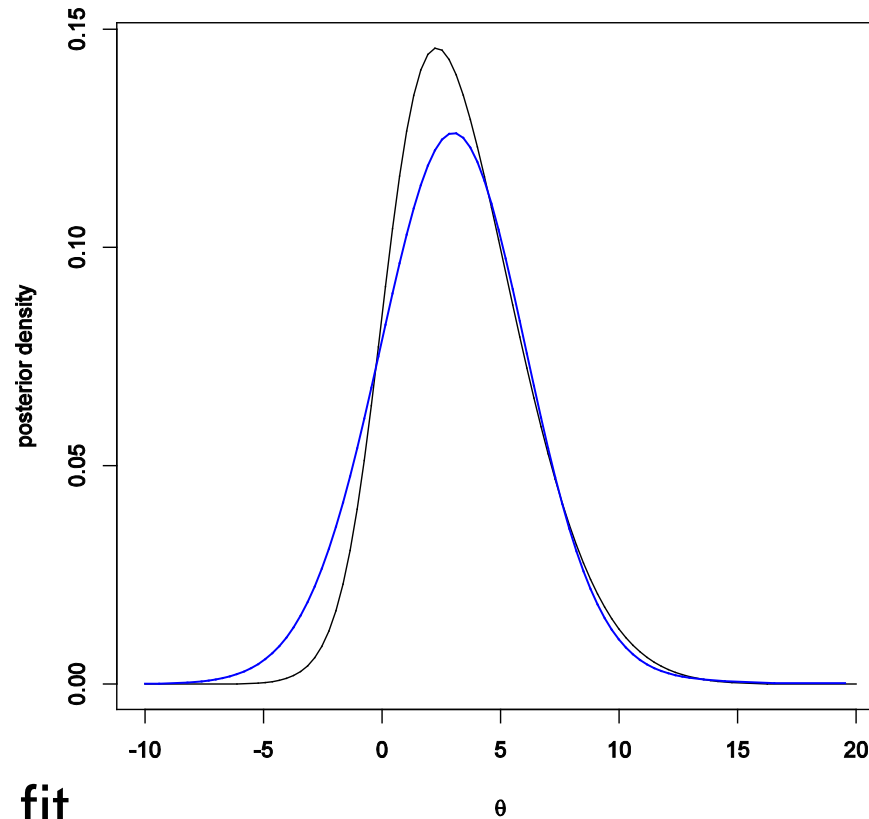
- Quadratic program.
- Then, DoIt

$$\hat{p}(\boldsymbol{\theta}|\mathbf{y}) \approx \frac{\hat{\mathbf{c}}' \boldsymbol{\phi}(\boldsymbol{\theta}; \boldsymbol{\Sigma})}{\hat{\mathbf{c}}' \mathbf{1}}$$

becomes a mixture normal approximation.

# Mixture Normal Approximation

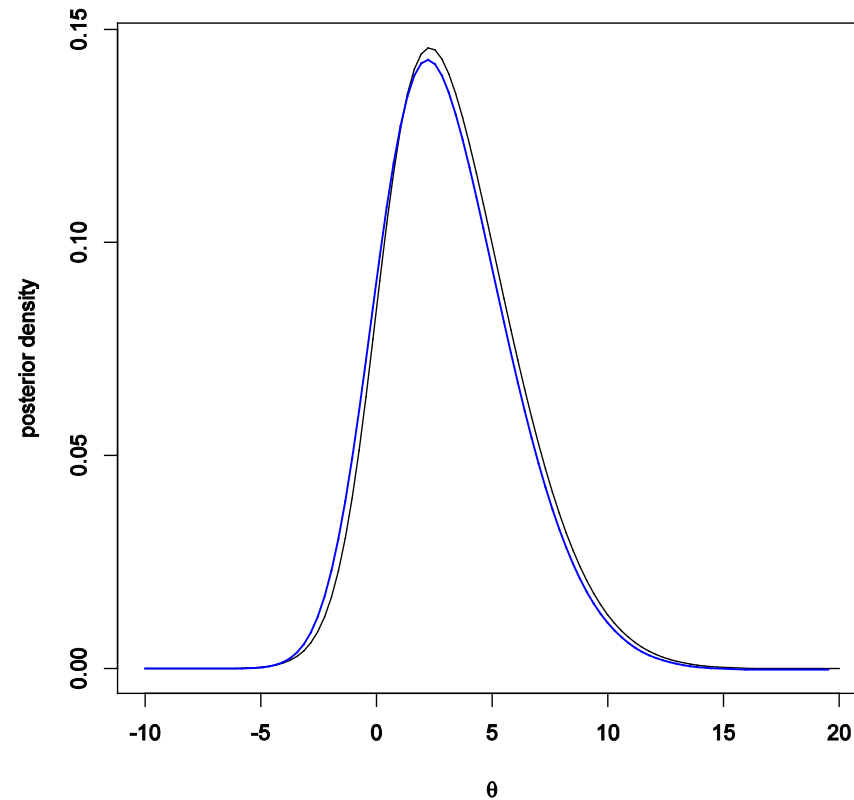
- $m=10$



- Not a good fit.

# Mixture Normal Approximation

- $m=20$



- Better, but not good enough.

# Improved DoIt

- DoIt :

$$h(\boldsymbol{\theta}) \approx \sum_{i=1}^m \hat{c}_i g(\boldsymbol{\theta}; \boldsymbol{\nu}_i, \boldsymbol{\Sigma}) \left\{ a + \sum_{i=1}^m b_i g(\boldsymbol{\theta}; \boldsymbol{\nu}_i, \boldsymbol{\Lambda}) \right\}$$

- Let  $z_i = h(\boldsymbol{\nu}_i) / \hat{\mathbf{c}}' \mathbf{g}(\boldsymbol{\nu}_i; \boldsymbol{\Sigma})$  for  $i = 1, \dots, m$ .

- Then,  $\hat{\mathbf{b}} = \mathbf{G}(\boldsymbol{\Lambda})^{-1}(\mathbf{z} - a\mathbf{1})$

- New approximation:

$$\hat{h}(\boldsymbol{\theta}) = \hat{\mathbf{c}}' \mathbf{g}(\boldsymbol{\theta}; \boldsymbol{\Sigma}) \{a + \hat{\mathbf{b}}' \mathbf{g}(\boldsymbol{\theta}; \boldsymbol{\Lambda})\}$$

# Improved DoIt-continued

$$\begin{aligned}\int \hat{h}(\boldsymbol{\theta}) d\boldsymbol{\theta} &= a \hat{\mathbf{c}}' \int g(\boldsymbol{\theta}; \boldsymbol{\Sigma}) d\boldsymbol{\theta} + \hat{\mathbf{c}}' \int g(\boldsymbol{\theta}; \boldsymbol{\Sigma}) g(\boldsymbol{\theta}; \boldsymbol{\Lambda})' d\boldsymbol{\theta} \hat{\mathbf{b}} \\ &= a \hat{\mathbf{c}}' (2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2} \mathbf{1} + \hat{\mathbf{c}}' (2\pi)^{d/2} \frac{|\boldsymbol{\Sigma} \boldsymbol{\Lambda}|^{1/2}}{|\boldsymbol{\Sigma} + \boldsymbol{\Lambda}|^{1/2}} \mathbf{G}(\boldsymbol{\Sigma} + \boldsymbol{\Lambda}) \hat{\mathbf{b}}\end{aligned}$$

$$a = \int \hat{z}(\boldsymbol{\theta}) \frac{\hat{\mathbf{c}}' \phi(\boldsymbol{\theta}; \boldsymbol{\Sigma})}{\hat{\mathbf{c}}' \mathbf{1}} d\boldsymbol{\theta} = \frac{\int \hat{h}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2} \hat{\mathbf{c}}' \mathbf{1}}$$

$$a = \frac{\hat{\mathbf{c}}' \mathbf{G}(\boldsymbol{\Sigma} + \boldsymbol{\Lambda}) \mathbf{G}(\boldsymbol{\Lambda})^{-1} \mathbf{z}}{\hat{\mathbf{c}}' \mathbf{G}(\boldsymbol{\Sigma} + \boldsymbol{\Lambda}) \mathbf{G}(\boldsymbol{\Lambda})^{-1} \mathbf{1}}$$

# Improved DoIt-continued

$$\int \hat{h}(\boldsymbol{\theta}) d\boldsymbol{\theta} = a(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2} \hat{\mathbf{c}}' \mathbf{1}$$

$$\hat{p}(\boldsymbol{\theta} | \mathbf{y}) \approx \frac{\hat{\mathbf{c}}' \phi(\boldsymbol{\theta}; \boldsymbol{\Sigma})}{\hat{\mathbf{c}}' \mathbf{1}} \{1 + \hat{\mathbf{b}}' \mathbf{g}(\boldsymbol{\theta}; \boldsymbol{\Lambda}) / a\}$$

- $\boldsymbol{\Lambda}$  can be obtained using cross validation.

# Improved DoIt-continued

Let  $V = \Sigma(\Sigma + \Lambda)^{-1}\Lambda$  and  $\mu_{ij} = V(\Sigma^{-1}\nu_i + \Lambda^{-1}\nu_j)$

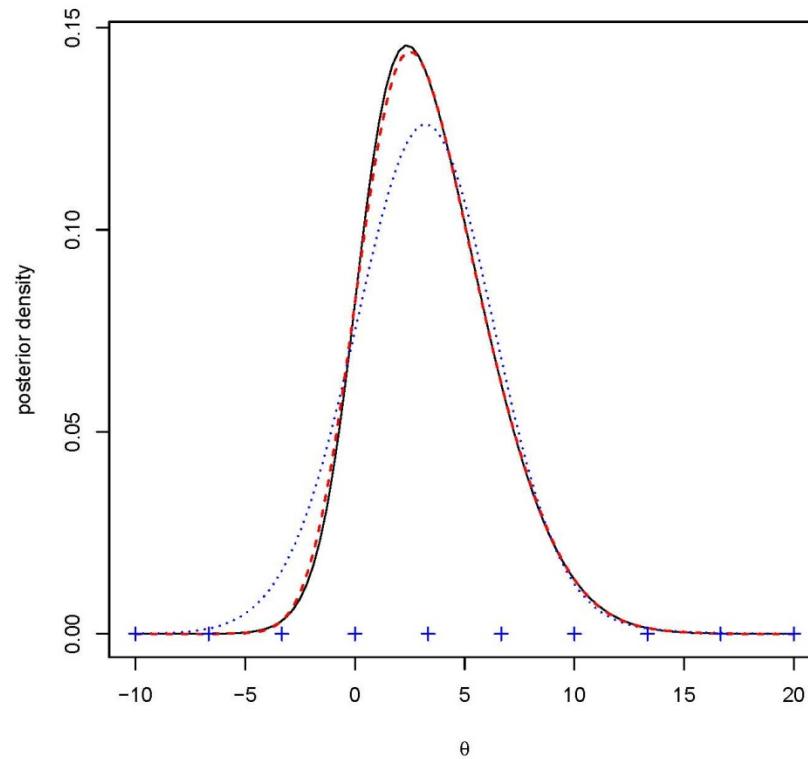
$$\hat{p}(\theta|y) \approx \frac{\sum_{i=1}^m \hat{c}_i \phi(\theta; \nu_i, \Sigma) + \sum_{i=1}^m \sum_{j=1}^m d_{ij} \phi(\theta; \mu_{ij}, V)}{\sum_{i=1}^m \hat{c}_i}$$

- where

$$d_{ij} = \frac{\hat{c}_i \hat{b}_j |\Lambda|^{1/2}}{a |\Sigma + \Lambda|^{1/2}} g(\nu_i; \nu_j, \Sigma + \Lambda).$$

# Binary Data Example

- $m=10$



# Another Example

- Marin and Robert (2007)

$$y|\theta \sim \text{Cauchy}(\theta, 1),$$

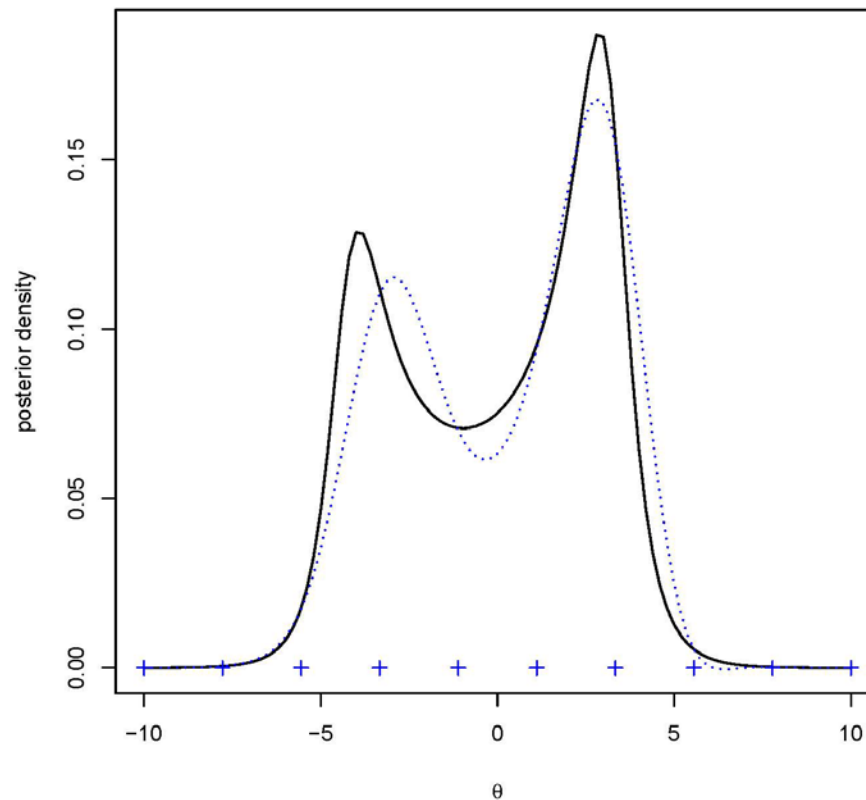
$$\theta \propto N(0, (\sqrt{10})^2).$$

- Unnormalized posterior:

$$h(\theta) = \frac{\exp(-\theta^2/20)}{\prod_{i=1}^2 (1 + (y_i - \theta)^2)}.$$

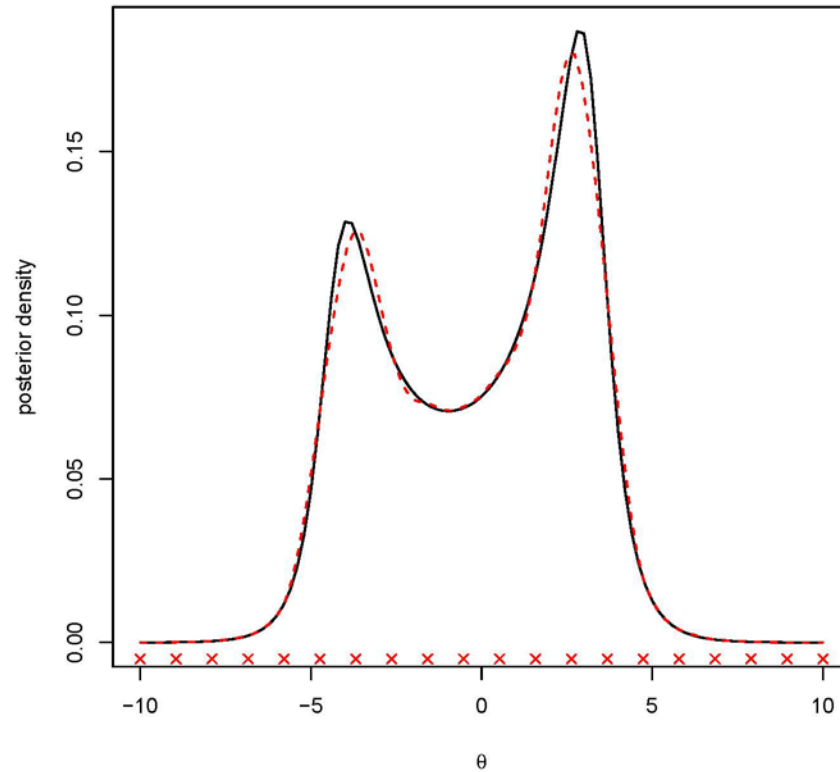
# Example-continued

Suppose we sample 10 equally spaced points from  $-10$  to  $10$ .



# Example-continued

- $m=20$



# Marginal Posterior Distributions

- Using properties of multivariate normal distribution:

$$\hat{p}(\theta_k | \mathbf{y}) \approx \frac{\sum_{i=1}^m \hat{c}_i \phi(\theta_k; \nu_{ik}, \Sigma_{kk}) + \sum_{i=1}^m \sum_{j=1}^m d_{ij} \phi(\theta_k; \mu_{ijk}, \mathbf{V}_{kk})}{\sum_{i=1}^m \hat{c}_i}$$

# Posterior Summaries

$$E(\boldsymbol{\theta}|\mathbf{y}) = \bar{\boldsymbol{\theta}} \approx \frac{\sum_{i=1}^m \hat{c}_i \boldsymbol{\nu}_i + \sum_{i=1}^m \sum_{j=1}^m d_{ij} \boldsymbol{\mu}_{ij}}{\sum_{i=1}^m \hat{c}_i}$$

$$\text{var}(\boldsymbol{\theta}|\mathbf{y}) \approx \frac{\sum_{i=1}^m \hat{c}_i (\boldsymbol{\nu}_i \boldsymbol{\nu}_i' + \boldsymbol{\Sigma}) + \sum_{i=1}^m \sum_{j=1}^m d_{ij} (\boldsymbol{\mu}_{ij} \boldsymbol{\mu}_{ij}' + \mathbf{V})}{\sum_{i=1}^m \hat{c}_i} - \bar{\boldsymbol{\theta}} \bar{\boldsymbol{\theta}}'$$

# Posterior Expectation

- More generally,

$$\begin{aligned}\xi &= E\{f(\boldsymbol{\theta})|\mathbf{y}\} \\ &\approx \int f(\boldsymbol{\theta}) \frac{\hat{\mathbf{c}}' \phi(\boldsymbol{\theta}; \boldsymbol{\Sigma})}{a \hat{\mathbf{c}}' \mathbf{1}} \{a + \hat{\mathbf{b}}' \mathbf{g}(\boldsymbol{\theta}; \boldsymbol{\Lambda})\} d\boldsymbol{\theta}\end{aligned}$$

for some function  $f(\boldsymbol{\theta})$ .

- Use approximations.

# Posterior Expectation-continued

- First, let  $z(\theta) = a + \hat{\mathbf{b}}' \mathbf{g}(\theta; \Lambda)$  and  $f^*(\theta) = f(\theta)z(\theta)$ .

$$\xi \approx \frac{1}{a\hat{\mathbf{c}}' \mathbf{1}} \sum_{i=1}^m \hat{c}_i \int f^*(\theta) \phi(\theta; \nu_i, \Sigma) d\theta$$

- Let  $\mathbf{f}^* = (f^*(\nu_1), \dots, f^*(\nu_m))' = \mathbf{f} \odot \mathbf{z}$
- Approximate  $f^*(\theta)$  using a kriging predictor:

$$f^*(\theta) = \alpha z(\theta) + \mathbf{g}(\theta; \Omega)' \mathbf{G}(\Omega)^{-1} (\mathbf{f}^* - \alpha \mathbf{z})$$

# Posterior Expectation-continued

$$\xi \approx \alpha + \frac{|\Omega|^{1/2}}{a\hat{c}'\mathbf{1}|\Omega + \Sigma|^{1/2}}\hat{c}'G(\Omega + \Sigma)G(\Omega)^{-1}(f^* - \alpha z)$$

- Choose  $\alpha = \xi$  (Joseph 2006),

$$\xi \approx \frac{\hat{c}'G(\Omega + \Sigma)G(\Omega)^{-1}f^*}{\hat{c}'G(\Omega + \Sigma)G(\Omega)^{-1}z}$$

- Take  $\Omega = \Lambda$ ,

$$\xi \approx \frac{\hat{c}'G(\Sigma + \Lambda)G(\Lambda)^{-1}f^*}{\hat{c}'G(\Sigma + \Lambda)G(\Lambda)^{-1}z}$$

# Example

- Posterior predictive density in the binary data example:

$$y|\mathbf{y} \sim \text{Bernoulli}(\xi),$$

where

$$\xi = E(\{1 + \exp(-\theta)\}^{-1}|\mathbf{y}).$$

- Numerical integration:  $\xi = .8496$
- Kriging approximation:  $\xi \approx .8478$
- First order approximation:  $\xi \approx \{1 + \exp(-\hat{\theta})\}^{-1} = .914$

# Experimental Design

- Initial space-filling design
- Sequential design

# Space-Filling Design

- By Laplace's approximation

$$\boldsymbol{\theta}|\mathbf{y} \sim^a N(\hat{\boldsymbol{\theta}}, \boldsymbol{\Sigma})$$

- Transform:

$$\boldsymbol{\alpha} = \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})$$

- Then,

$$\boldsymbol{\alpha}|\mathbf{y} \sim^a N(\mathbf{0}, \mathbf{I})$$

First choose a design  $\mathbf{D}^* = (\boldsymbol{\nu}_1^*, \dots, \boldsymbol{\nu}_m^*)'$  from  $(0, 1)^d$

$$\mathbf{D} = (\hat{\boldsymbol{\theta}} + \boldsymbol{\Sigma}^{1/2}\Phi^{-1}(\boldsymbol{\nu}_1^*), \dots, \hat{\boldsymbol{\theta}} + \boldsymbol{\Sigma}^{1/2}\Phi^{-1}(\boldsymbol{\nu}_m^*))'$$

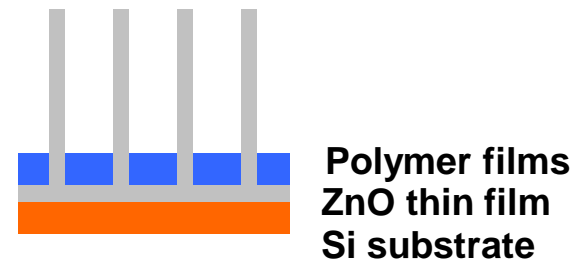
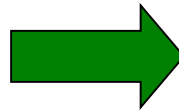
# Space-Filling Design-continued

- Latin Hypercube Design (LHD)
- Maximin LHD (Morris and Mitchell 1995)
- Let  $\nu_1^* = .5 = (.5, \dots, .5)'$
- Find the remaining  $m-1$  points by minimizing

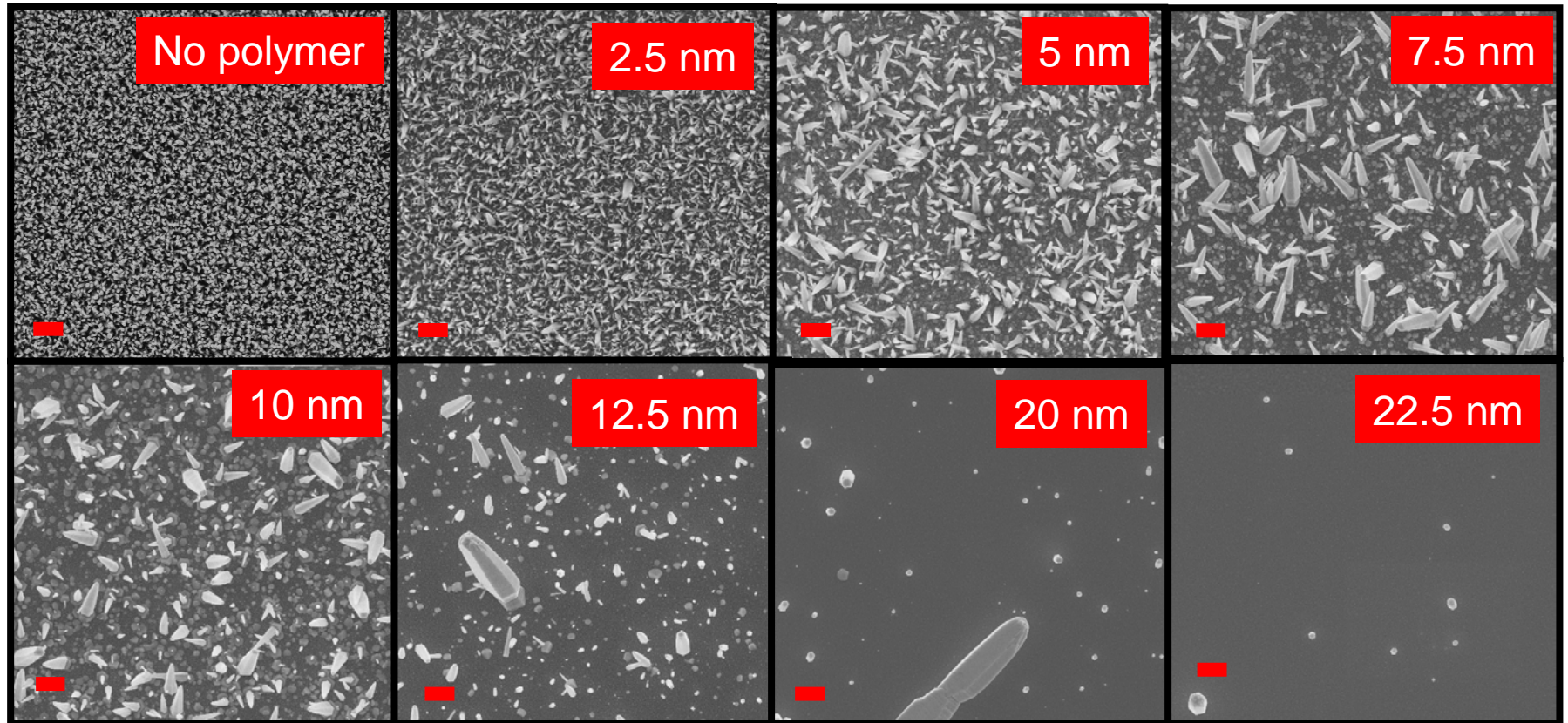
$$\left\{ \sum_{i=2}^m \sum_{j=2}^m 1/d^k(\nu_i^*, \nu_j^*) \right\}^{1/k}$$

# Example: Density Control of Nanowires

- Dasgupta, Weintraub, and Joseph (2011)



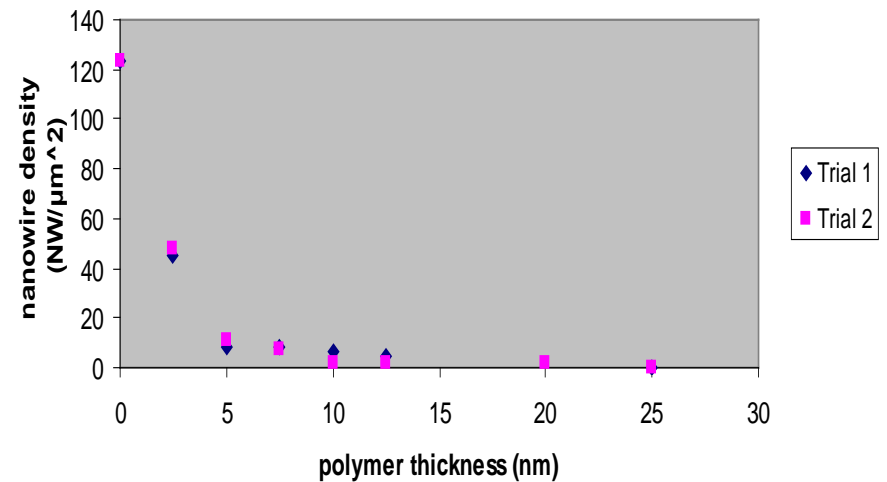
# Results (1st set of experiments)



(All images at 5000x magnification;  $\text{—}$  = 1  $\mu\text{m}$ )

# Experimental Data

# Bi Layers	Thickn ess (nm)	Density (NW/ $\mu\text{m}^2$ )	
		Trial 1	Trial 2
0	0.0	123	123
1	2.5	46	48
2	5.0	8	11
3	7.5	8	7
4	10.0	7	2
5	12.5	5	1
8	20.0	2	-
10	25.0	0	0



## Example-continued

- Density of nanowires ( $y$ )
- Thickness of polymer films ( $x$ )

$$y_{ij} | \boldsymbol{\theta}, \mathbf{u} \sim \text{Poisson}(\mu(x_i))$$

$$\mu(x_i) = [\theta_1 \exp(-\theta_2 x_i^2) + \theta_3 \{1 - \exp(-\theta_2 x_i^2)\} \Phi(-x_i/\theta_4)] u_i$$

for  $i = 1, \dots, 8$  and  $j = 1, 2$ .

- Let  $\gamma_i = \log(\theta_i)$  for  $i = 1, \dots, 4$ .  
and  $\alpha_i = \log(u_i)$  for  $i = 1, \dots, 8$
- Prior:  $p(\boldsymbol{\gamma}) \propto 1$ , and  $\alpha_i \sim^{iid} N(0, .1^2)$

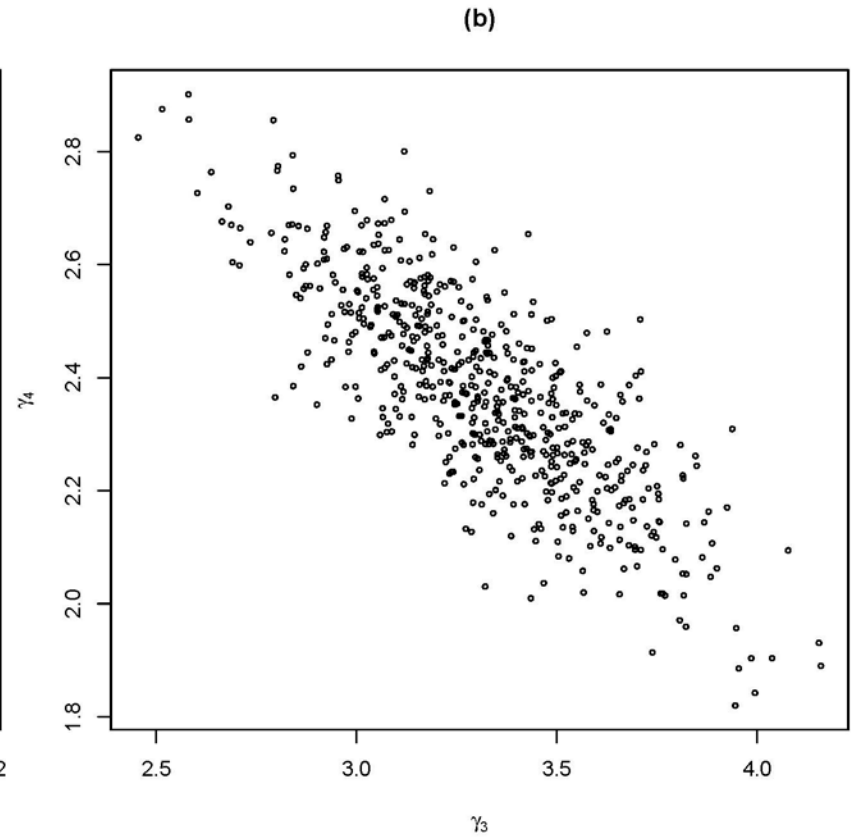
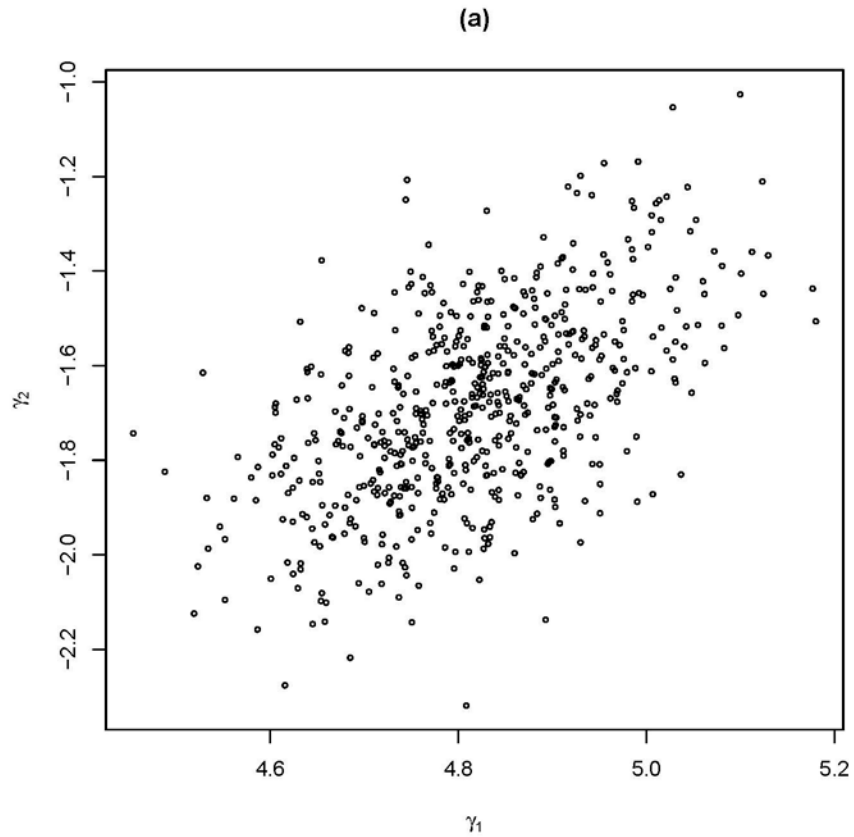
## Example-continued

- Posterior mode:  $\hat{\gamma} = (4.82, -1.69, 3.32, 2.37)'$

$$\hat{\alpha} = (-0.003, 0.005, -0.008, 0.014, -0.007, -0.007, 0.011, -0.005)'$$

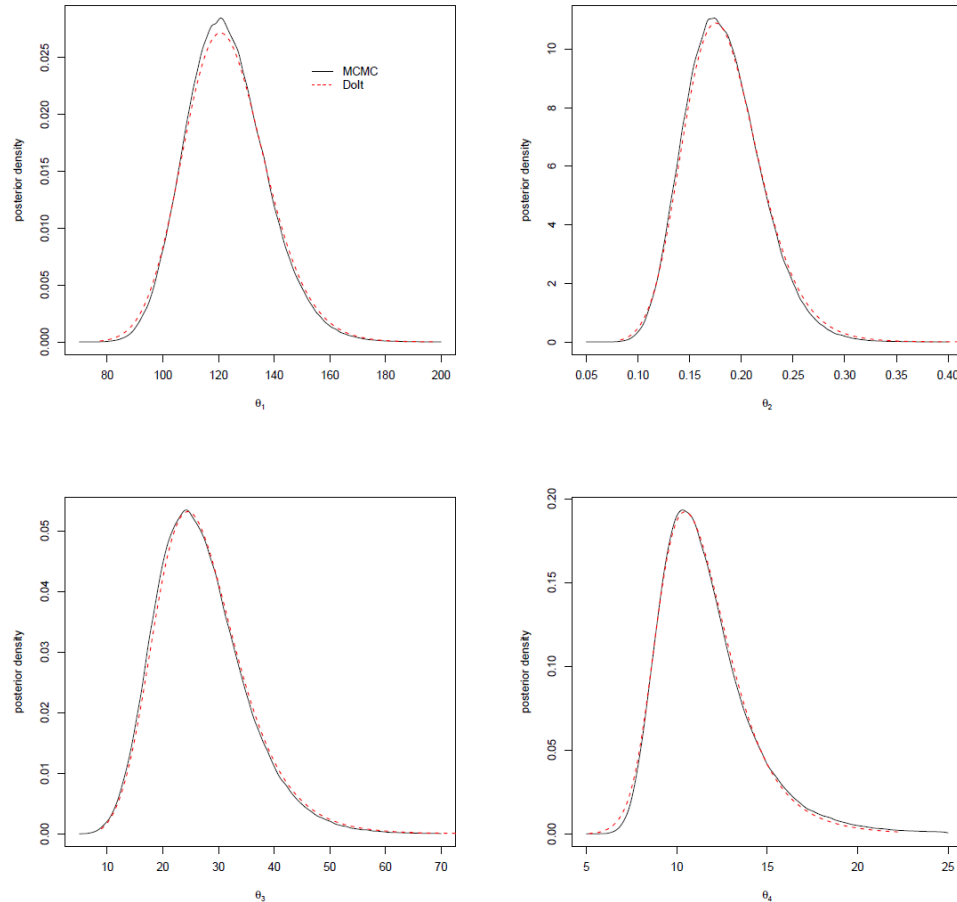
- $\Sigma$  obtained through numerical differentiation.
- $m = 50 \times 12 = 600$
- MmLHD
- $D = (\hat{\theta} + \Sigma^{1/2}\Phi^{-1}(\nu_1^*), \dots, \hat{\theta} + \Sigma^{1/2}\Phi^{-1}(\nu_m^*))'$

# Example-continued



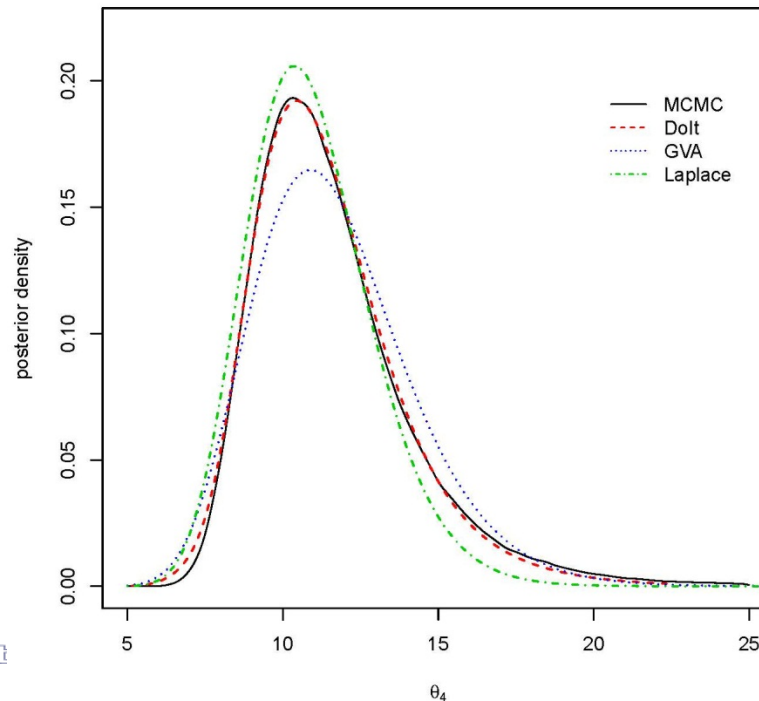
# Example-continued

- Metropolis algorithm with 3,000,000 samples (black line).



# Comparison

- Quadrature method: didn't converge!
- Gaussian Variational Approximation (GVA) (Ormerod and Wand 2012)



# Sequential Design

- Add points one-by-one to improve approximation.
- Optimal design theory: add the new point at a location with the largest prediction uncertainty.
- DoIt can be viewed as a simple kriging predictor given  $\hat{c}'g(\theta; \Sigma)$ .

- Conditional prediction variance is

$$\left(\hat{c}'g(\theta; \Sigma)\right)^2 \{1 - g(\theta; \Lambda)'G^{-1}(\Lambda)g(\theta; \Lambda)\}$$

$$\nu_{m+1} = \arg \max_{\theta} \left(\hat{c}'g(\theta; \Sigma)\right)^2 \{1 - g(\theta; \Lambda)'G^{-1}(\Lambda)g(\theta; \Lambda)\}$$

# Sequential Design-continued

- Let  $v_{(i)}$  be the leave-one-out estimate of the prediction variance

$$v_{(i)} \approx \left( h_i + l_i - \frac{G_i^{-1}(\Sigma)}{G_{ii}^{-1}(\Sigma)}(h + l) \right)^2 \frac{1}{G_{ii}^{-1}(\Lambda)}$$

for  $i = 1, \dots, m$ , where  $l = G(\Sigma)\hat{c} - h$

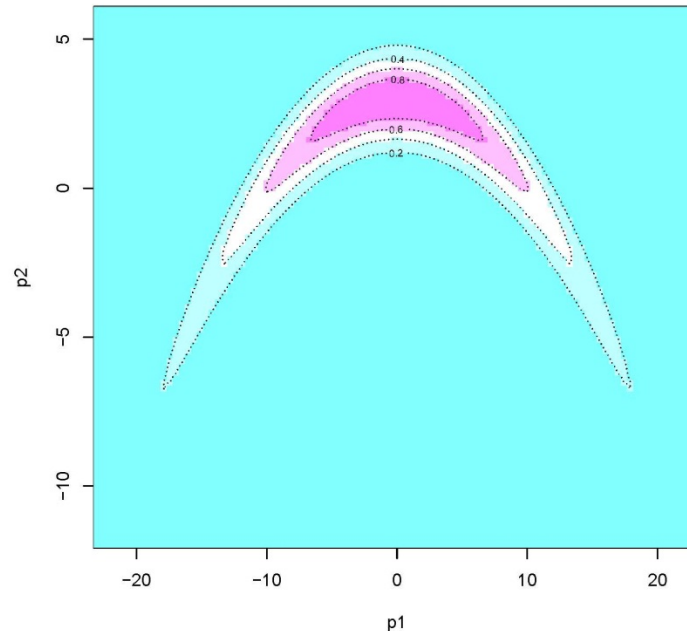
Let  $i^* = \arg \max_i v_{(i)}$ .

- Optimize in the neighborhood of  $\nu_{i^*}$

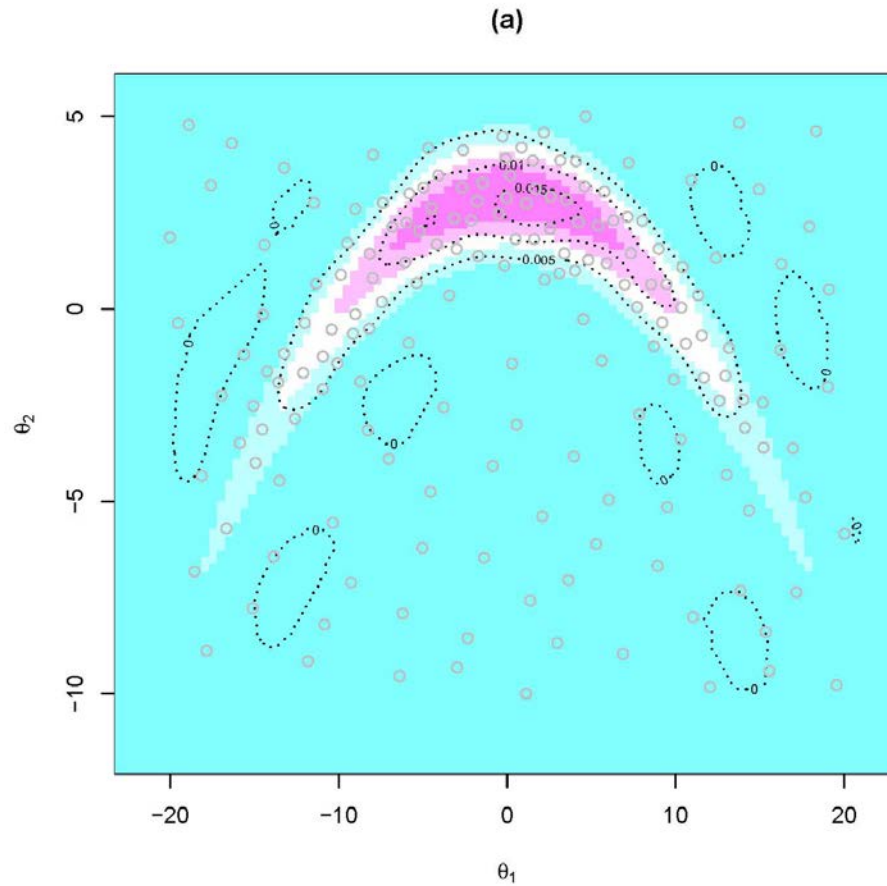
# Example

- Haario, Saksman, and Tamminen (2001)

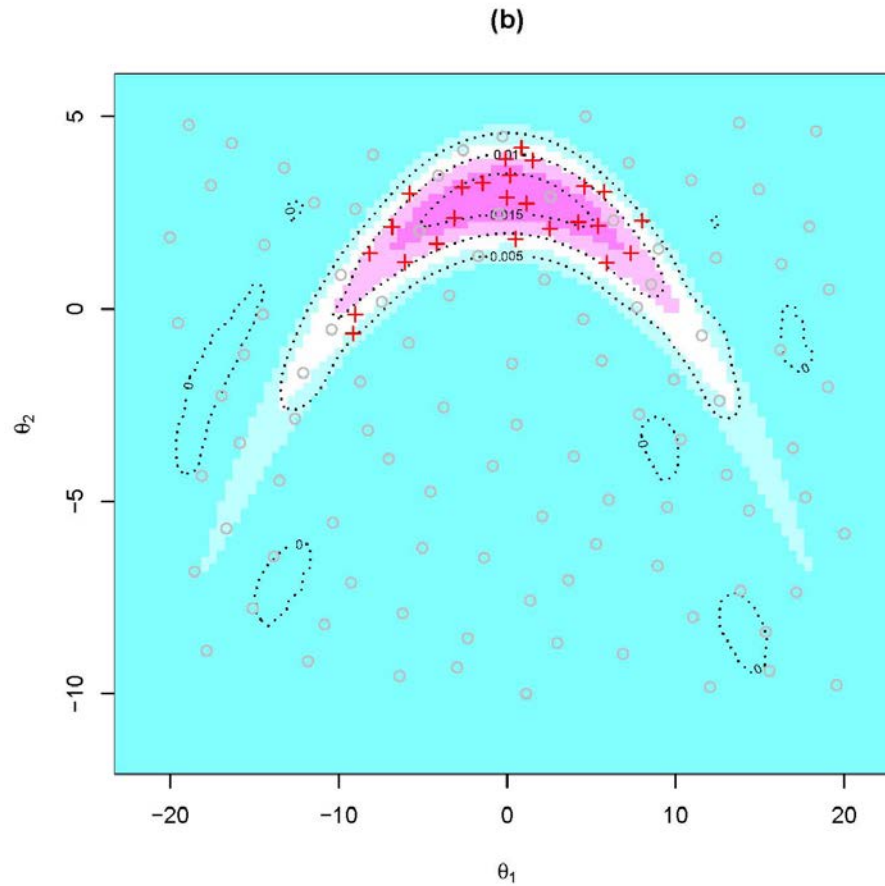
$$p(\boldsymbol{\theta}|\mathbf{y}) = \phi \left( (\theta_1, \theta_2 + .03\theta_1^2 - 3)'; (0, 0)', \text{diag}\{100, 1\} \right)$$



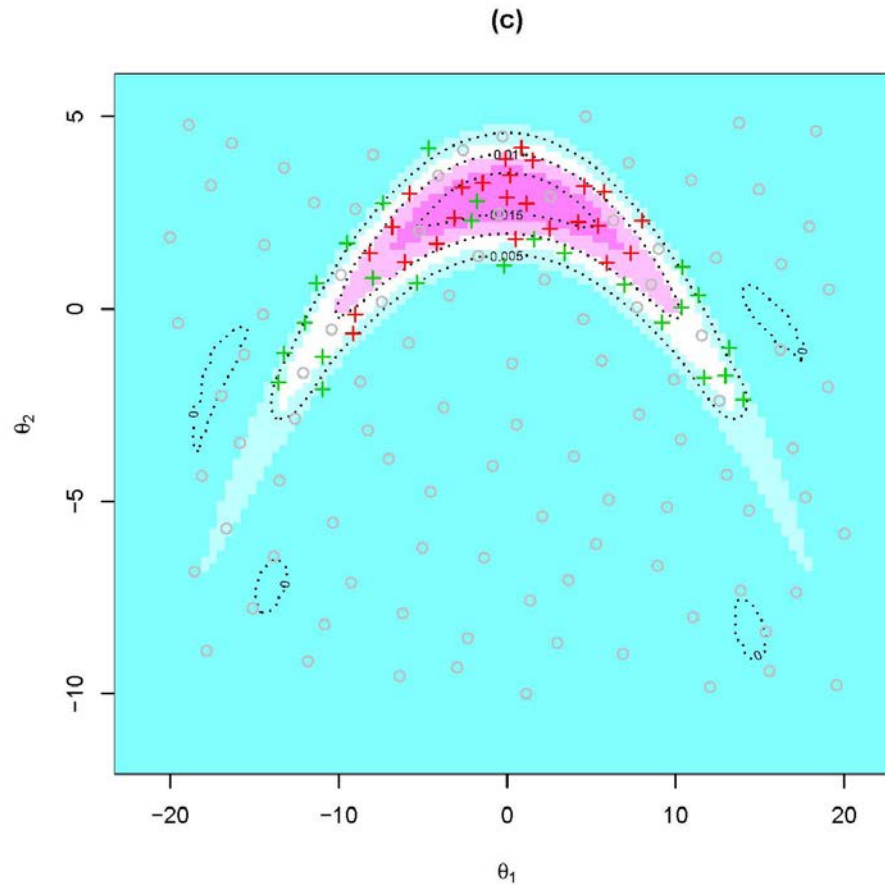
$m=100$



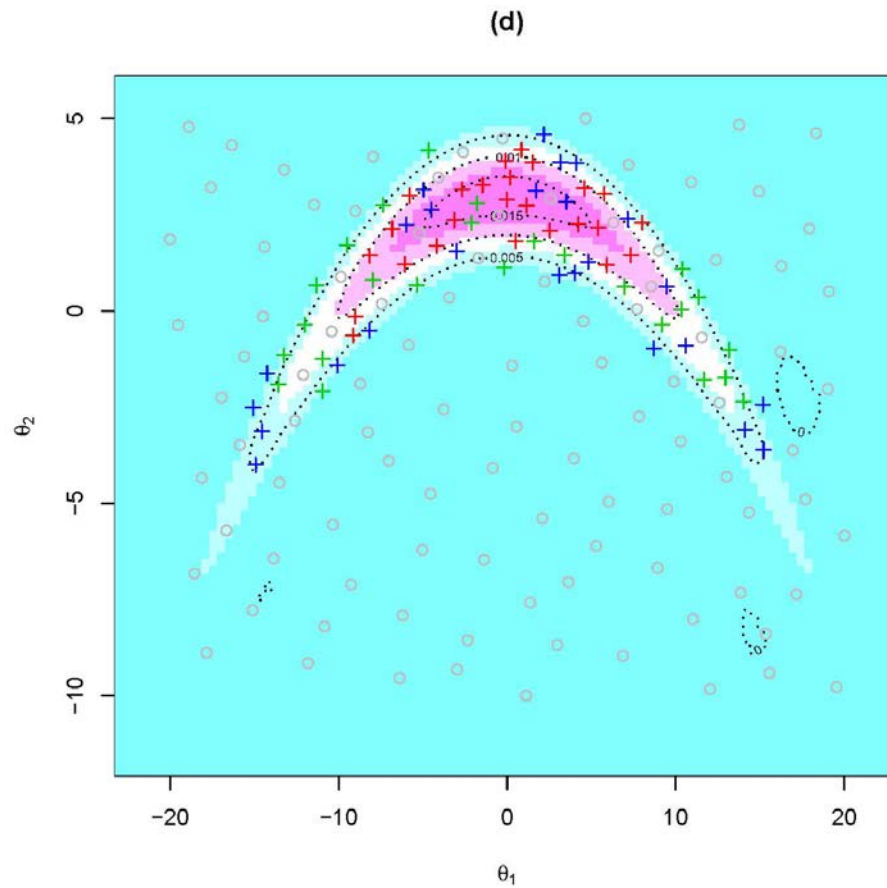
$$m=100+25$$



$$m=100+50$$



$$m=100+75$$

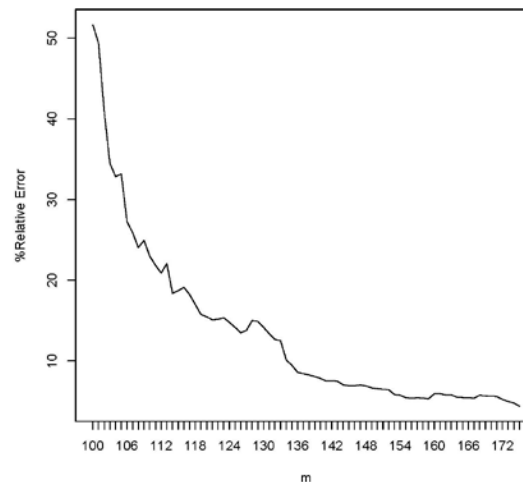


# Example-continued

- % Relative Error:  $\%RE = \frac{\overline{|cv|}}{\bar{h}} \times 100,$

where  $\overline{|cv|} = E(|cv(\boldsymbol{\theta})| | \mathbf{y})$  and  $\bar{h} = E(h(\boldsymbol{\theta}) | \mathbf{y})$

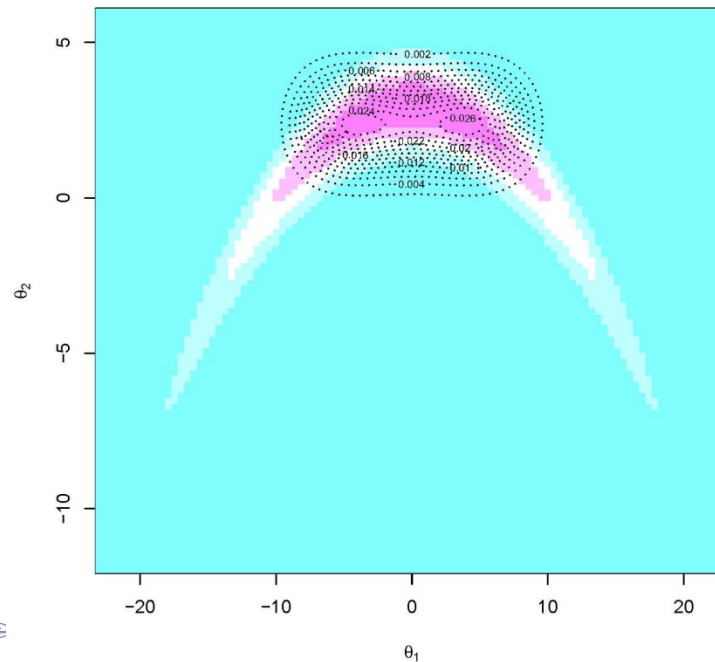
$$cv_i = h_i - \left( h_i + l_i - \frac{G_i^{-1}(\Sigma)}{G_{ii}^{-1}(\Sigma)}(h + l) \right) \left( \frac{h_i}{h_i + l_i} - \frac{G_i^{-1}(\Lambda)}{G_{ii}^{-1}(\Lambda)} \left( \frac{h}{h + l} - a1 \right) \right)$$



# Comparison: Variational Bayes

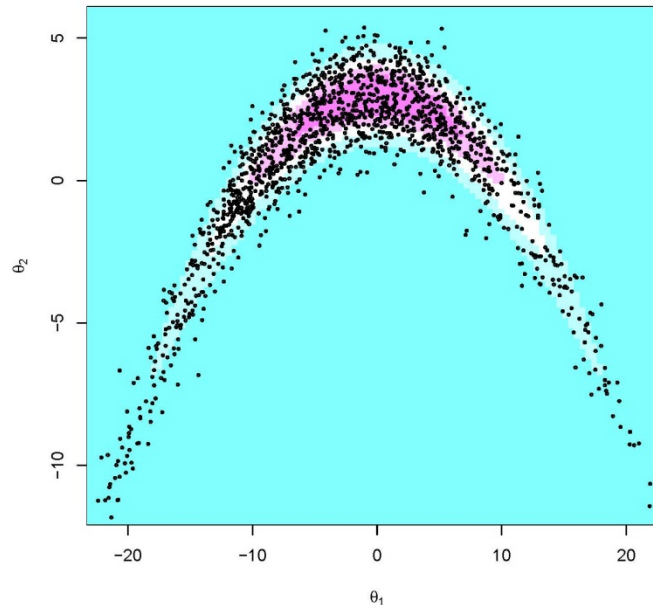
$$\hat{p}_{VB}(\boldsymbol{\theta}|\mathbf{y}) = q_1(\theta_1)\phi(\theta_2; \mu_2, \sigma_2^2),$$

where  $q_1(\theta_1) \propto \exp\{-.5[\theta_1^2/100 + (\mu_2 + .03\theta_1^2 - 3)^2]\}$ ,  
 $\mu_2 = -.03(\mu_1^2 + \sigma_1^2) + 3$ , and  $\sigma_2^2 = 1$ .



# Comparison: Hybrid MCMC

- Fielding, Nott, and Liong (2011)



- CPU time: Hybrid MCMC=90 mins, Dolt=3 mins.

## Discussion from Dagupta and Meng: connections to QMC

- Experimental design
  - QMC: low discrepancy sequence in  $[0,1]^d$
  - Dolt: initial space-filling design+ sequential design
- Posterior summaries
  - QMC: Monte Carlo average
  - Dolt: smooth interpolation + analytical evaluation
- Dolt is likely to perform better than QMC when the posterior densities are smooth.

# Numerical Comparison

- Binary data example

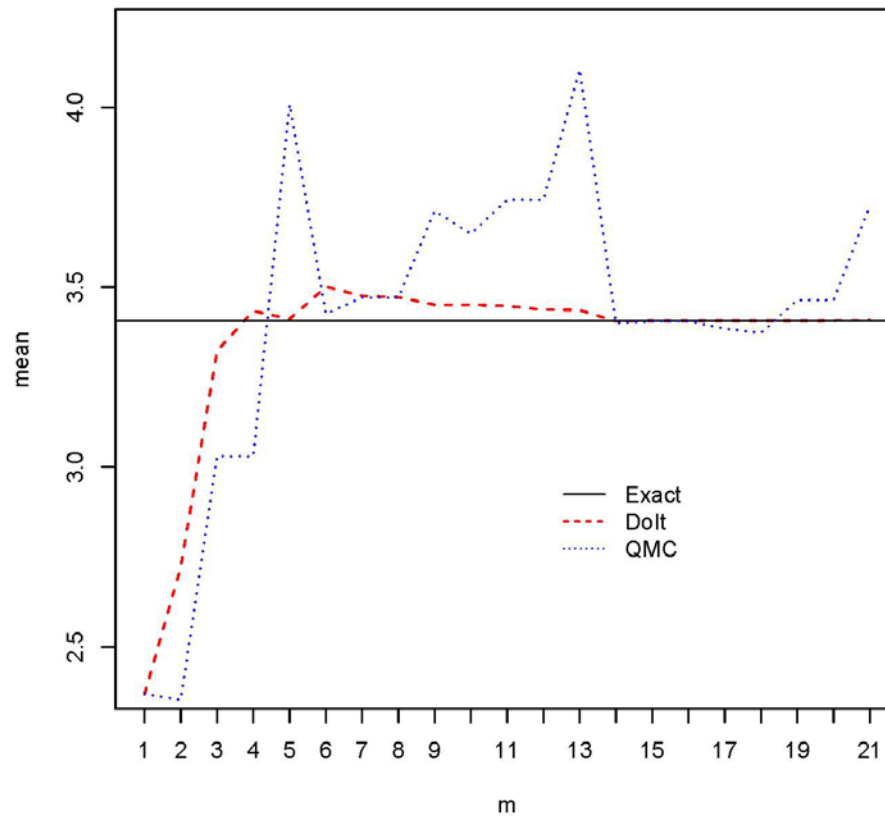
$$y|\theta \sim \text{Bernoulli}(\{1 + \exp(-\theta)\}^{-1}),$$

$$\theta \sim N(\mu, \tau^2).$$

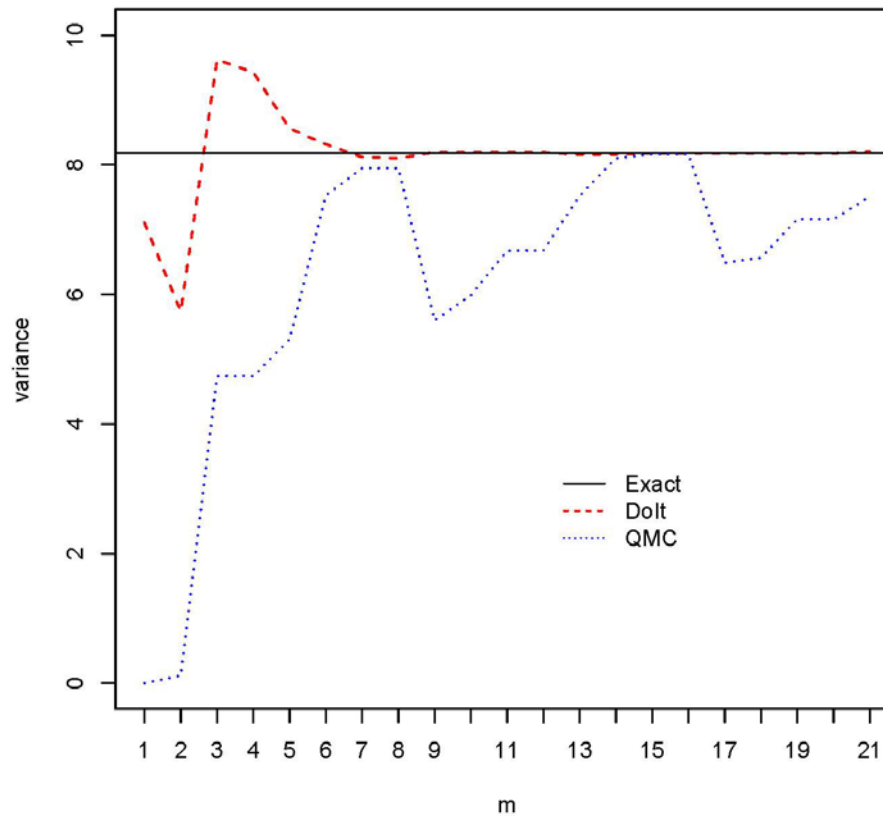
- Laplace approximation:  $N(2.37, 2.67^2)$
- van der Corput sequence:  $\frac{1}{2}, \frac{1}{4}, \frac{3}{4}, \dots$  and re-scaled to  $[2.37-15, 2.37+15]$ .

$$\hat{\theta}_{QMC} = E(\theta|y)_{QMC} = \frac{\sum_{i=1}^m h_i \nu_i}{\sum_{i=1}^m h_i} \text{ and } \text{var}(\theta|y)_{QMC} = \frac{\sum_{i=1}^m h_i (\nu_i - \hat{\theta}_{QMC})^2}{\sum_{i=1}^m h_i}.$$

# Posterior Mean



# Posterior Variance



# Hierarchical Models

- May contain very large number of parameters
- Not easy to find a good space-filling design in high dimensions.
  - Make use of the hierarchical structure.

$$y|\theta \sim p(y|\theta), \theta|\eta \sim p(\theta|\eta), \text{ and } \eta \sim p(\eta).$$

# Hierarchical Models-continued

- Suppose we can obtain explicit expression of

$$p(\mathbf{y}|\boldsymbol{\eta}) = \int p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\eta}) d\boldsymbol{\theta},$$

- And have the conditional distribution

$$p(\boldsymbol{\theta}|\boldsymbol{\eta}, \mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\eta})$$

- Use DoIt :  $\hat{p}(\boldsymbol{\eta}|\mathbf{y}) \approx \frac{\hat{\mathbf{c}}' \phi(\boldsymbol{\eta}; \boldsymbol{\Sigma})}{\hat{\mathbf{c}}' \mathbf{1}} \{1 + \hat{\mathbf{b}}' \mathbf{g}(\boldsymbol{\eta}; \boldsymbol{\Lambda})/a\}$ .

- Then,  $\hat{p}(\boldsymbol{\theta}|\mathbf{y}) \approx \frac{\hat{\mathbf{c}}' \mathbf{G}(\boldsymbol{\Sigma} + \boldsymbol{\Lambda}) \mathbf{G}(\boldsymbol{\Lambda})^{-1} \mathbf{p}^*(\boldsymbol{\theta})}{\hat{\mathbf{c}}' \mathbf{G}(\boldsymbol{\Sigma} + \boldsymbol{\Lambda}) \mathbf{G}(\boldsymbol{\Lambda})^{-1} \mathbf{z}}$

where  $\mathbf{p}^*(\boldsymbol{\theta}) = (p(\boldsymbol{\theta}|\boldsymbol{\nu}_1, \mathbf{y}), \dots, p(\boldsymbol{\theta}|\boldsymbol{\nu}_m, \mathbf{y}))' \odot \mathbf{z}$

# Example 1: A Longitudinal Data Analysis

- Orthodontic measurements on 27 children (Phinheiro and Bates 2000)

$$y_{ij} | \beta, u_i, \sigma_\epsilon^2 \sim^{ind.} N(\beta_0 + u_i + \beta_1 age_{ij} + \beta_2 sex_i, \sigma_\epsilon^2),$$

$$u_i | \sigma_u^2 \sim^{iid} N(0, \sigma_u^2),$$

$$\beta \sim N(\mathbf{0}, 10^8 \mathbf{I}_3),$$

$$\sigma_\epsilon^2, \sigma_u^2 \propto^{ind.} IG(.01, .01),$$

for  $i = 1, \dots, 27$  and  $j = 1, \dots, 4$ .

- 32 parameters**

# Longitudinal Data-continued

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$$

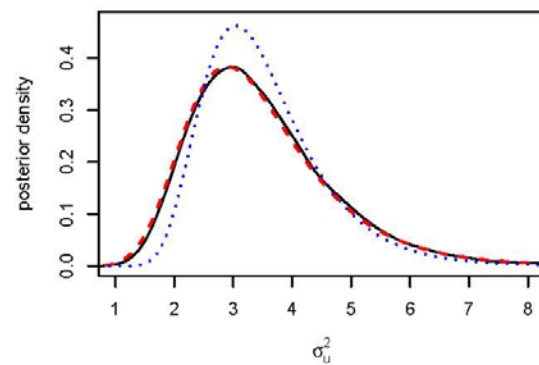
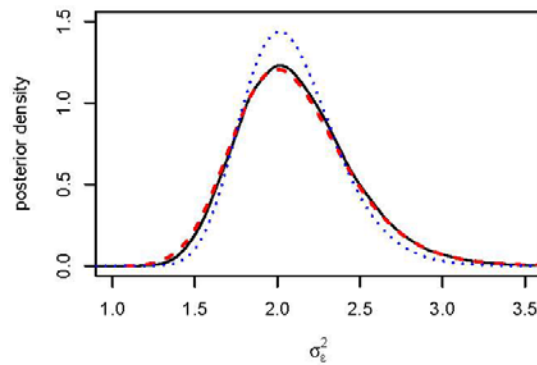
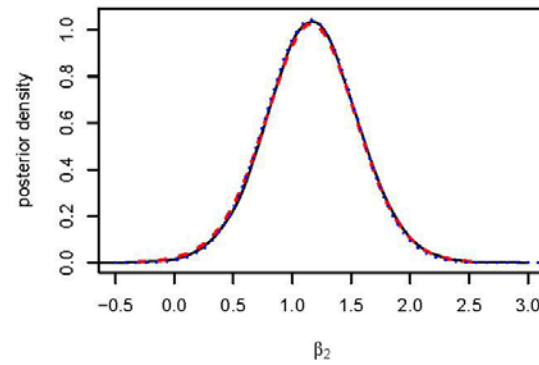
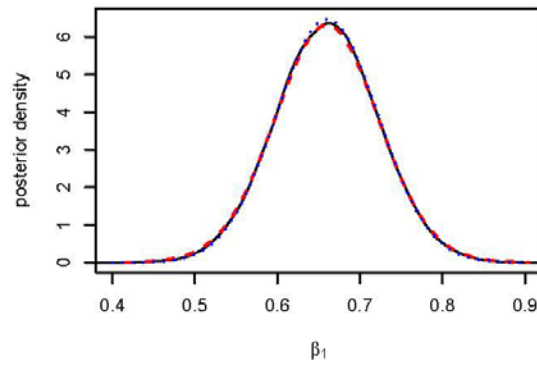
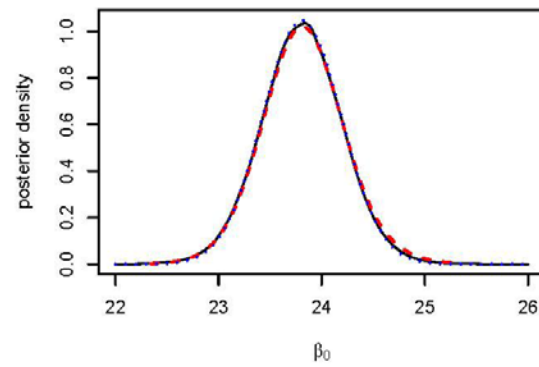
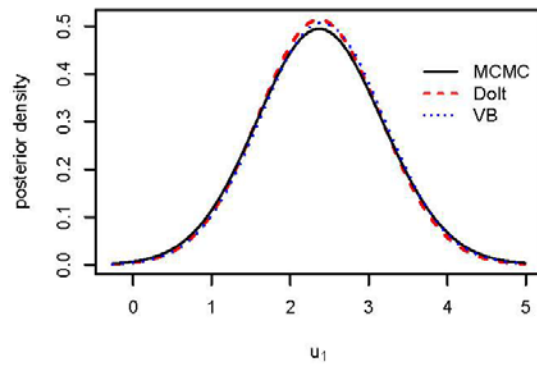
- Integrating out  $\mathbf{u}$

$$\mathbf{y}|\boldsymbol{\beta}, \sigma_{\epsilon}^2, \sigma_u^2 \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma_{\epsilon}^2 \mathbf{I}_{108} + \sigma_u^2 \mathbf{Z}\mathbf{Z}')$$

- Also,

$$\mathbf{u}|\boldsymbol{\beta}, \sigma_{\epsilon}^2, \sigma_u^2, \mathbf{y} \sim N((\mathbf{Z}'\mathbf{Z} + \frac{\sigma_{\epsilon}^2}{\sigma_u^2} \mathbf{I}_{27})^{-1} \mathbf{Z}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}), (\mathbf{Z}'\mathbf{Z} + \frac{\sigma_{\epsilon}^2}{\sigma_u^2} \mathbf{I}_{27})^{-1})$$

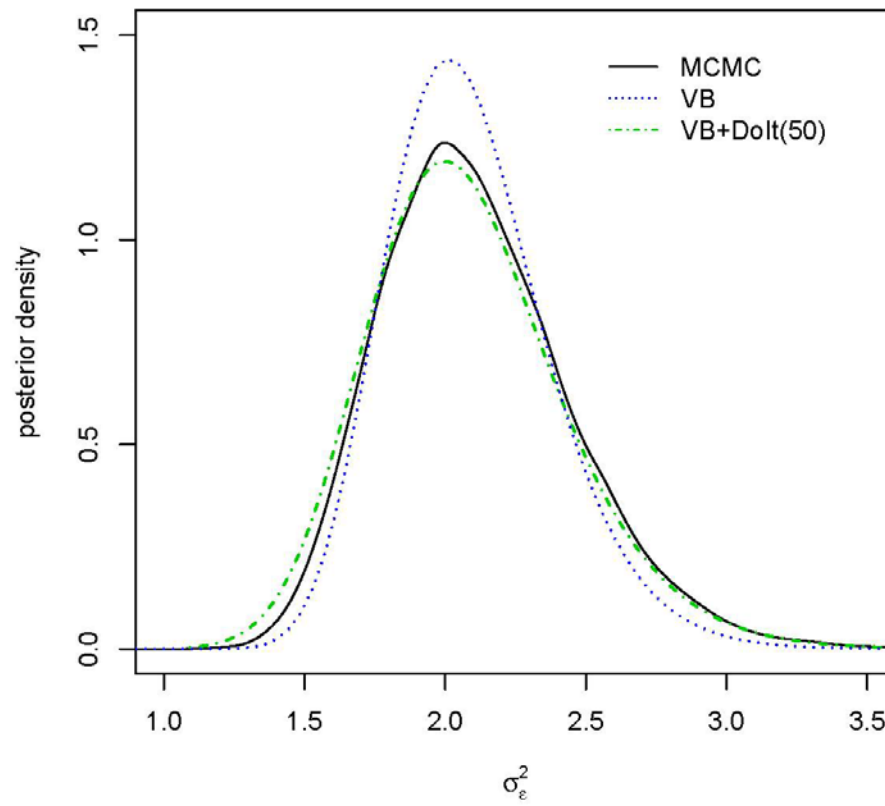
- DoIt : 250-run MmLHD in **5 dimensions**



# Discussion from Ormerod & Wand

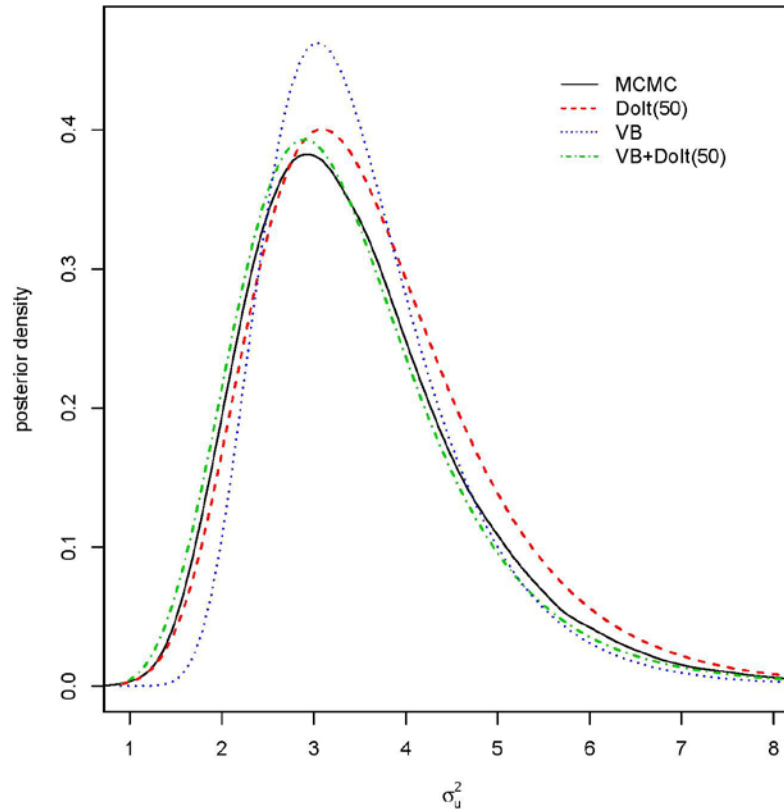
- VB approximation can be improved using the grid-based variational approximation method in Ormerod (2011).
- Dolt can be used for the same purpose!
  - Center the experimental design using VB estimates.

# VB+DoIt



# VB+Dolt

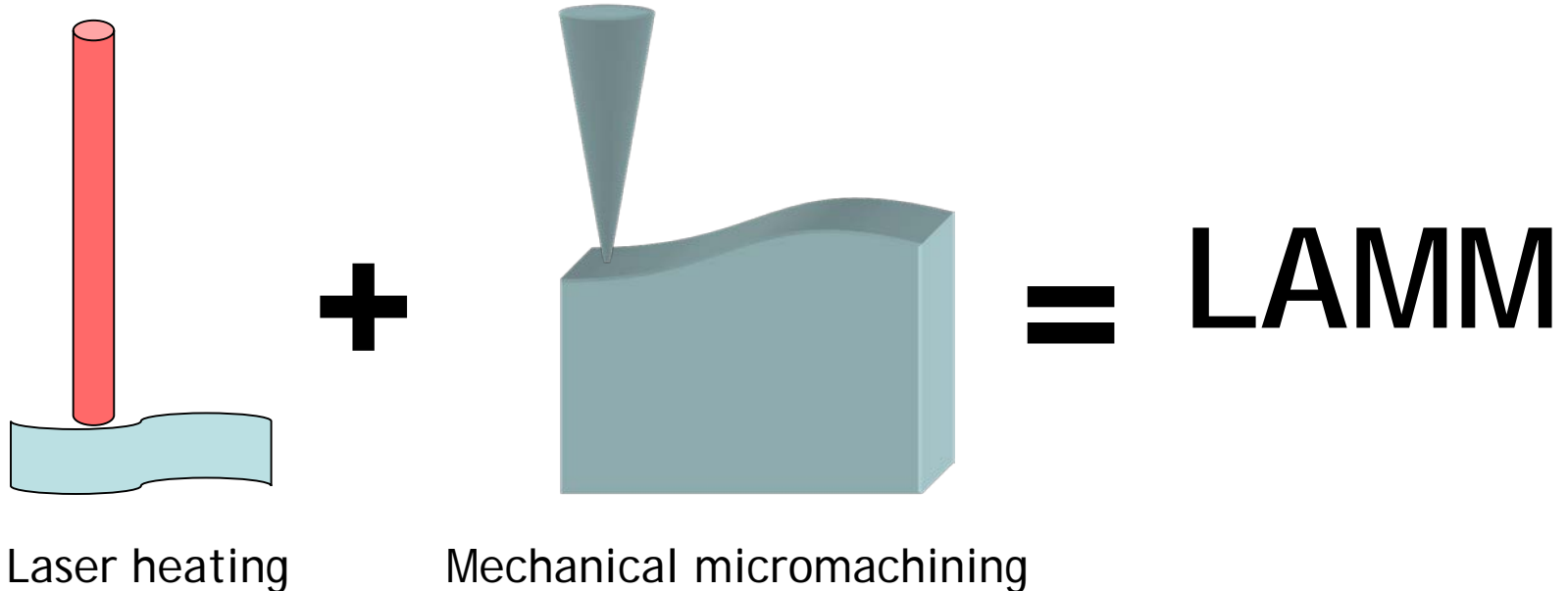
- VB can be used to improve Dolt, whenever VB implementation is readily available.



## Example 2: LAMM

Laser assisted mechanical micromachining (LAMM) integrates *thermal softening* with *mechanical micro cutting*

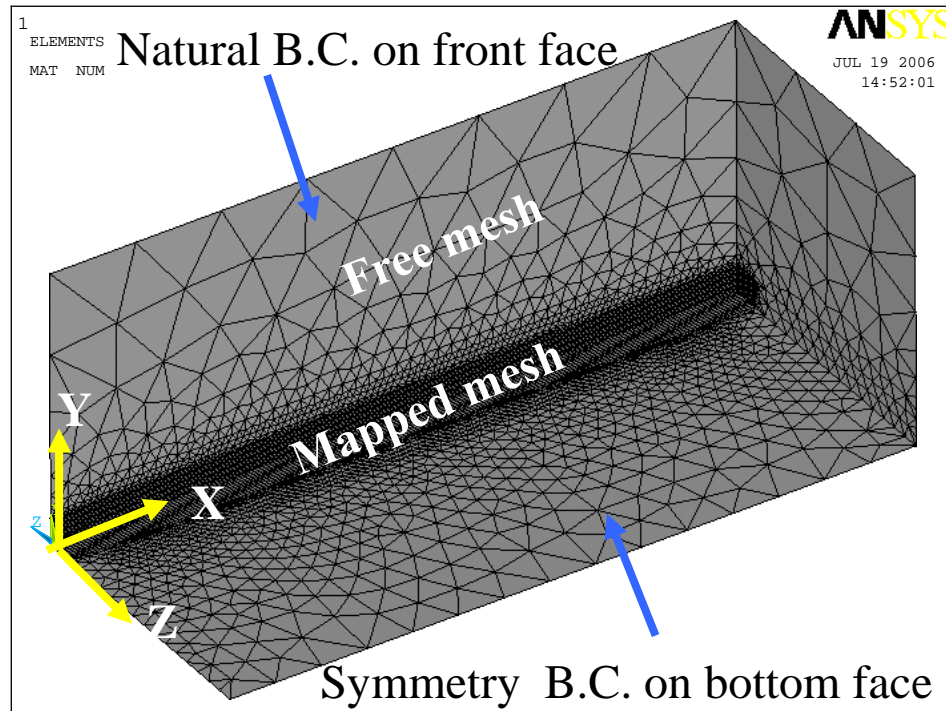
(Singh, Joseph, and Melkote 2011)



# Objective

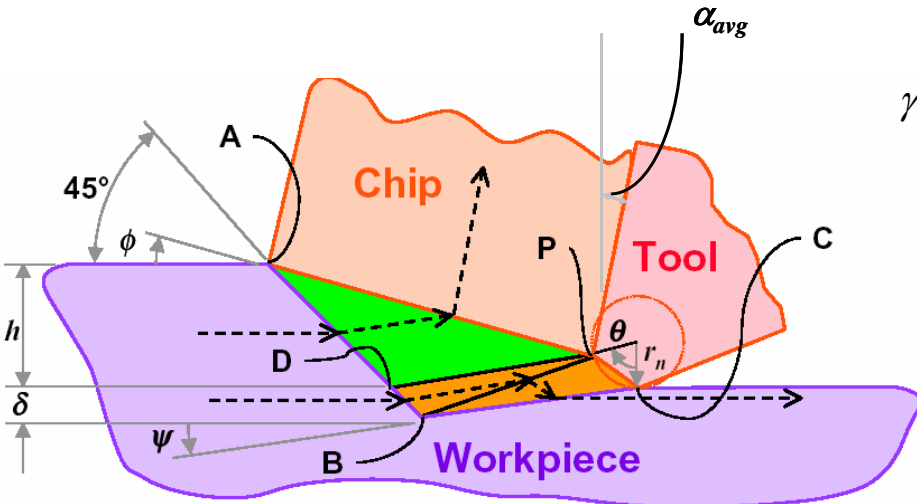
Find optimum processing conditions that minimize cutting/thrust forces and thermal damage.

# Thermal Model



- Mapped dense mesh ( $25\ \mu\text{m} \times 12.5\ \mu\text{m} \times 20\ \mu\text{m}$ )
- An 8 noded 3-D thermal element (Solid70) is used
- Gaussian distribution of heat flux applied to a 5x5 element matrix which sweeps the mesh on the front face

# Geometric Model



(Manjunathiah et. al, 2000)

$$\dot{\gamma}_{chip} = 2V \frac{\gamma_{chip}}{\sqrt{2} \sin(\pi/4 + \theta_{PD}) PD}$$

$$\dot{\gamma}_{work} = 2V \frac{\gamma_{work}}{\sqrt{2} \sin(\pi/4 + \theta_{PD}) \overline{PD} + \frac{\sin(\psi + \theta/2)}{\sin \psi} \overline{PC}}$$

$$\gamma_{chip} = \frac{\sqrt{2} \sin \theta_{PD}}{\sin(\pi/4 + \theta_{PD})} + \frac{\cos(\alpha_{avg} + \theta_{PD})}{\cos(\alpha_{avg} - \phi) \sin(\phi + \theta_{PD})}$$

$$\gamma_{work} = \frac{\sqrt{2} \sin \theta_{PD}}{\sin(\pi/4 + \theta_{PD})} + \frac{\sin(\theta_{PD} + \theta/2)}{\sin(\theta_{PB} + \theta/2) \sin(\theta_{PB} + \theta_{PD})} + \frac{\sin \theta/2}{\sin \psi \sin(\psi + \theta/2)}$$

$$\gamma_{eff} = \frac{v_{chip} \gamma_{chip} + v_{work} \gamma_{work}}{v_{chip} + v_{work}}$$

$$\gamma_{eff}^{\bullet} = \frac{v_{chip} \gamma_{chip} + v_{work} \gamma_{work}}{v_{chip} + v_{work}}$$

For plane strain conditions,

$$\varepsilon = \gamma_{eff} / \sqrt{3}$$

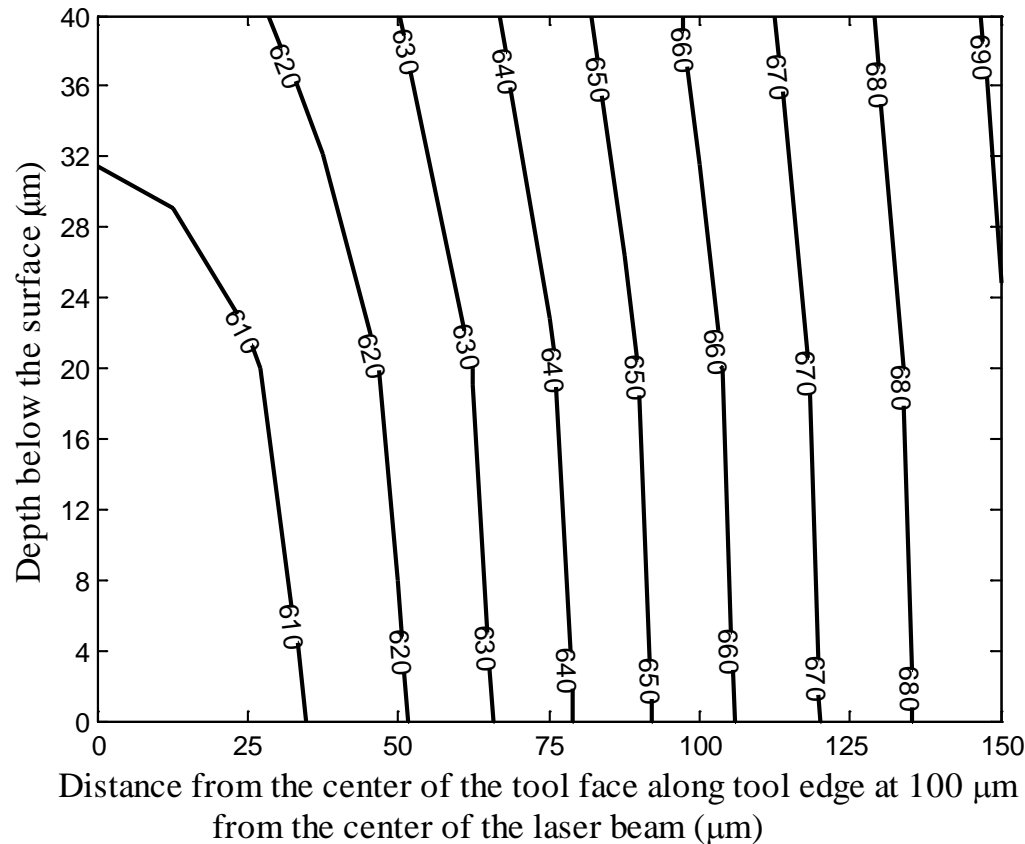
$$\dot{\varepsilon} = \dot{\gamma}_{eff} / \sqrt{3}$$

# Shear Flow Strength

$$\sigma(\varepsilon, \dot{\varepsilon}, T, HRC) = \left( A + B\varepsilon^n + C \ln(\varepsilon + \varepsilon_0) + D \right) \left( 1 + E \ln \left( \frac{\dot{\varepsilon}}{\dot{\varepsilon}_0} \right) \right) \left( 1 - (T^*)^m \right)$$

Yan et al., 2007

$$S = \sigma / \sqrt{3}$$



10W laser power, 10 mm/min speed 100 μm laser-tool distance  
and 110 μm spot size

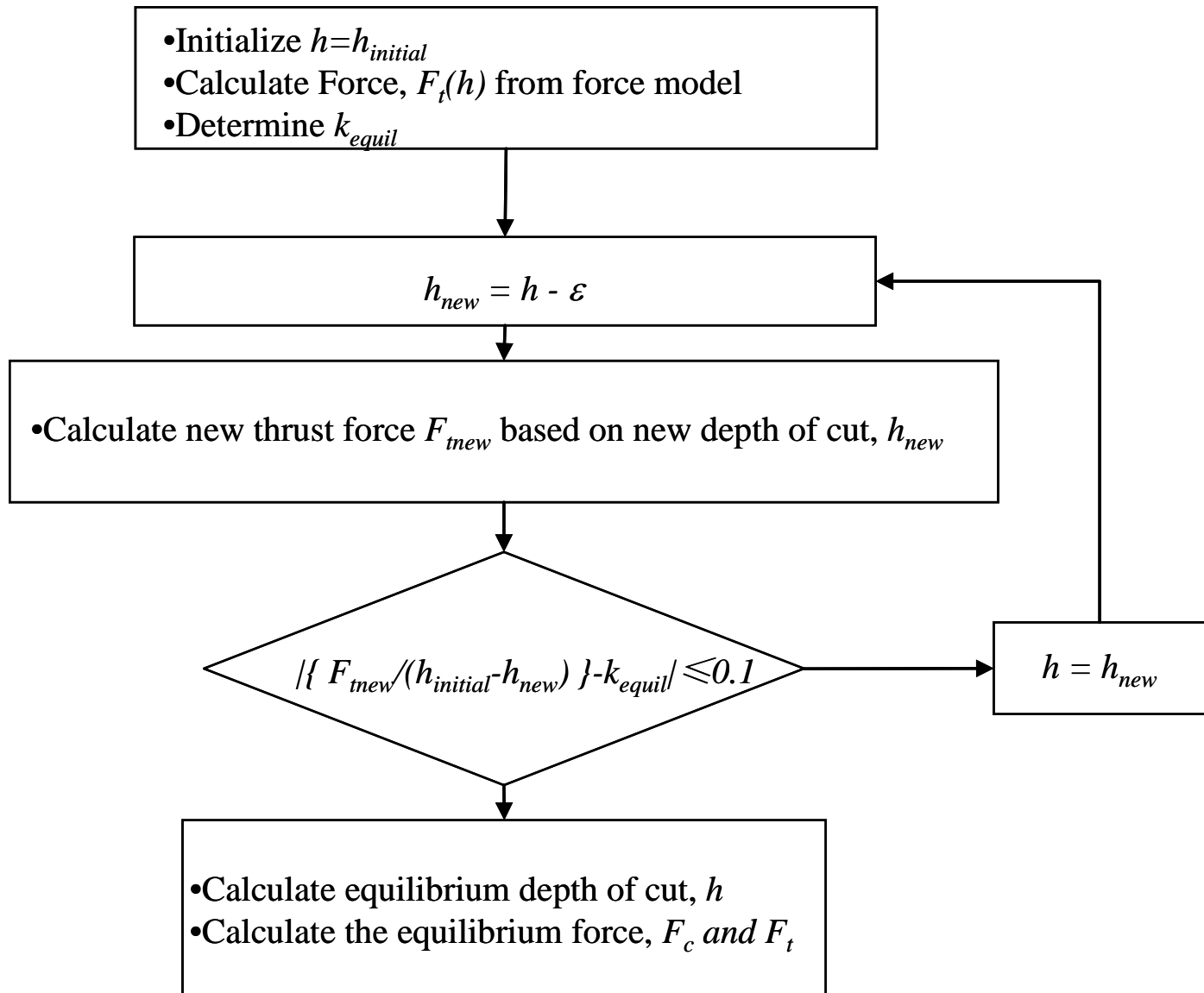
# Forces

- Cutting and thrust forces,

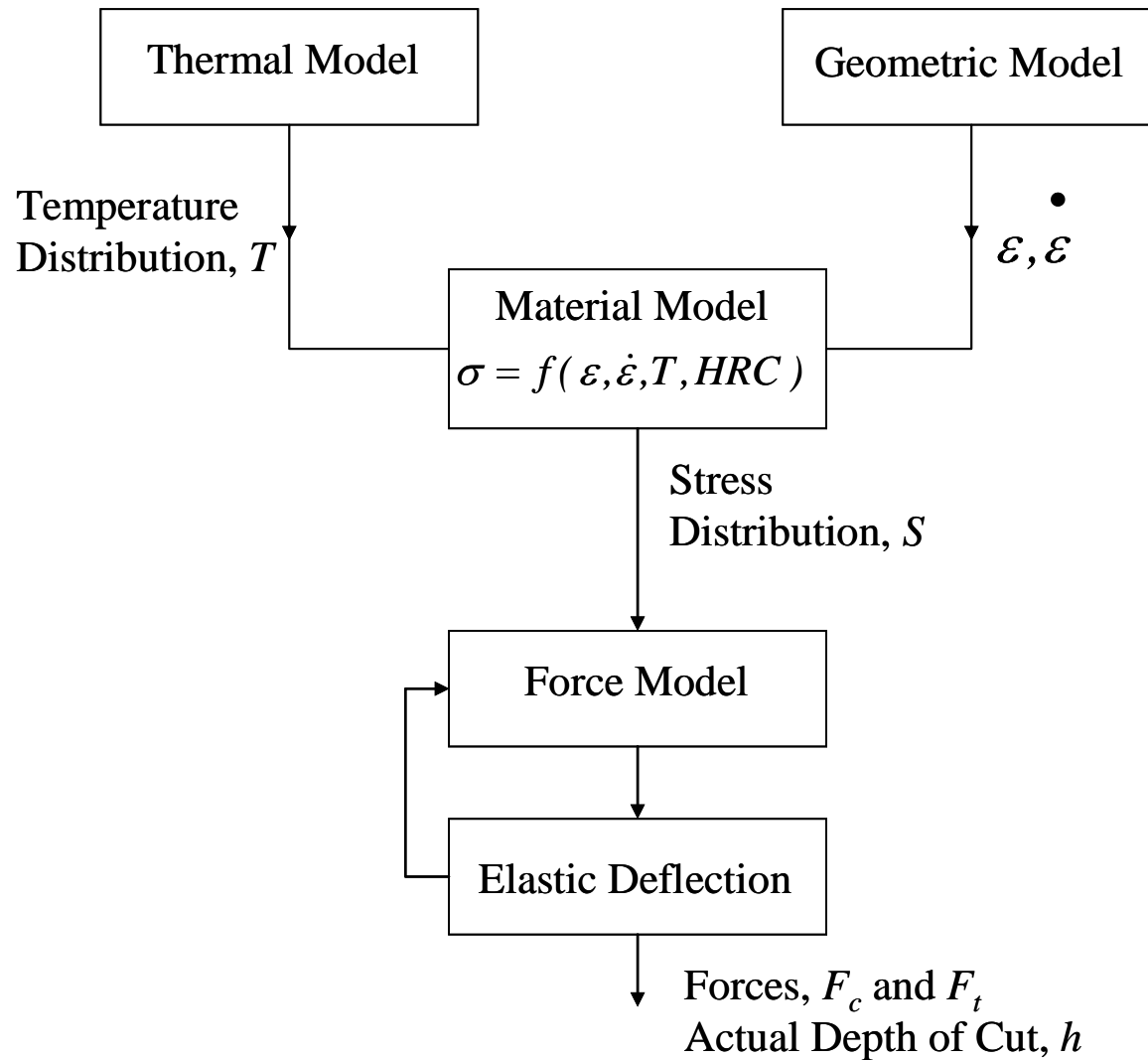
$$F_c = \{ (h - p) \cot \phi + h + r_n \sin \theta - (k - 1) \delta \} \sum_{i=1}^n \bar{S}(i) w(i)$$

$$F_t = \{ (h - p) \cot \phi - h + r_n \sin \theta + (k - 1) \delta \cot \psi \} \sum_{i=1}^n \bar{S}(i) w(i)$$

# Equilibrium Forces/Deflection



# Force model



# LAMM-continued

$$y_i = \theta(\mathbf{x}_i)$$

$$\theta(\mathbf{x}) \sim GP(\mu, \tau^2 r), \quad r(\mathbf{x}_i, \mathbf{x}_j) = \exp\left\{-\sum_{k=1}^4 \alpha_k (x_{ik} - x_{jk})^2\right\}$$

$$p(\mu, \tau^2) \propto 1/\tau^2$$

$$\gamma_i = \log(\alpha_i) \sim^{iid} N(0, 1).$$

# LAMM-continued

- Unnormalized posterior

$$\begin{aligned}h(\theta(\mathbf{x}), \mu, \tau^2, \gamma) &\propto p(\mathbf{y}|\theta(\mathbf{x}), \mu, \tau^2, \gamma)p(\theta(\mathbf{x})|\mu, \tau^2, \gamma)p(\mu, \tau^2)p(\gamma), \\ &= p(\theta(\mathbf{x})|\mu, \tau^2, \gamma, \mathbf{y})p(\mathbf{y}|\mu, \tau^2, \gamma)p(\mu, \tau^2)p(\gamma)\end{aligned}$$

- Integrating out  $\theta(\mathbf{x})$ ,  $\mu$ , and  $\tau^2$ ,

$$h(\gamma) = |\mathbf{R}|^{-1/2}(\mathbf{1}'\mathbf{R}^{-1}\mathbf{1})^{-1/2}[(\mathbf{y} - \hat{\mu}\mathbf{1})'\mathbf{R}^{-1}(\mathbf{y} - \hat{\mu}\mathbf{1})]^{-(n-1)/2}$$

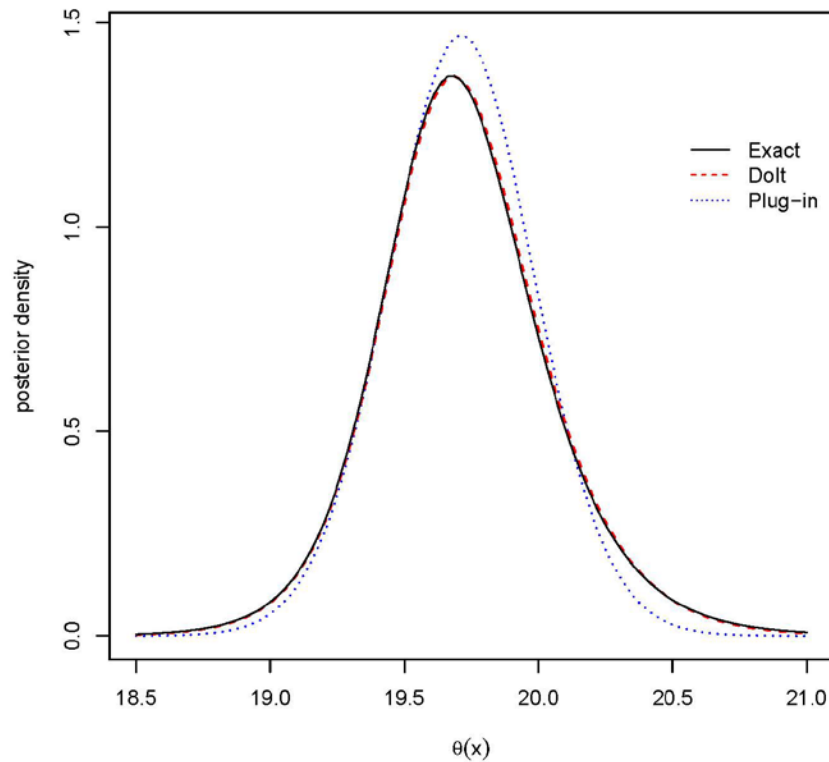
# LAMM-continued

- Conditional distribution:  $\frac{\theta(x) - \hat{\theta}(x)}{\sqrt{V(x)}} | \gamma, \mathbf{y} \sim t_{n-1},$

$$\begin{aligned}\hat{\theta}(x) &= \hat{\mu} + \mathbf{r}(x)' \mathbf{R}^{-1} (\mathbf{y} - \hat{\mu} \mathbf{1}), \\ V(x) &= \hat{\tau}^2 \left( 1 - \mathbf{r}(x)' \mathbf{R}^{-1} \mathbf{r}(x) + \frac{\{1 - \mathbf{r}(x)' \mathbf{R}^{-1} \mathbf{1}\}^2}{\mathbf{1}' \mathbf{R}^{-1} \mathbf{1}} \right), \\ \hat{\mu} &= \frac{\mathbf{1}' \mathbf{R}^{-1} \mathbf{y}}{\mathbf{1}' \mathbf{R}^{-1} \mathbf{1}}, \\ \hat{\tau}^2 &= \frac{1}{n-1} (\mathbf{y} - \hat{\mu} \mathbf{1})' \mathbf{R}^{-1} (\mathbf{y} - \hat{\mu} \mathbf{1}).\end{aligned}$$

# LAMM-continued

- DoIt :  $m=100$  (red), MCMC:  $m=100,000$  (black)



# LAMM-continued

- Computational time
  - DoIt = 3 seconds
  - MCMC=10 minutes



200 times faster!

- Computationally expensive likelihood.
  - Computational complexity:  $O(n^3)$ .
  - Here  $n=48$ .

# Conclusions

- A new deterministic approximation method using **design of experiments** and **interpolation** techniques.
  - Very general.
  - Can obtain the results with arbitrary precision.
  - Suffer less from the curse of dimensionality compared to lattice-based quadrature methods.
  - Some of the very large hierarchical problems can be solved efficiently.
  - Very fast!
  - Almost a black box method (no extra programming or derivations are needed).

# Conclusions

- Disadvantages
  - Not as flexible as MCMC.
  - Not easy to find good space-filling designs in high dimensions.
  - Can handle only continuous parameters.
  - Cannot ensure nonnegativity of the posterior density.

# Nonnegative Dolt

- Joseph, *Technometrics*, 2013, February.

Let  $\hat{\mathbf{c}} = \mathbf{b}$ ,  $\Lambda = \Sigma$ , and  $a = 0$ .

$$h(\boldsymbol{\theta}) \approx \left\{ \sum_{i=1}^m b_i g(\boldsymbol{\theta}; \boldsymbol{\nu}_i, \Sigma) \right\}^2,$$

which will always be nonnegative!

$$\sqrt{h(\boldsymbol{\theta})} \approx \sum_{i=1}^m b_i g(\boldsymbol{\theta}; \boldsymbol{\nu}_i, \Sigma)$$

# Future Research

- Fast generation of design points at high probability regions.
  - Need a better design strategy.
- Fast approximation using local versions of  $\Sigma$  .
  - Need a better modeling strategy.
- Topics for future research!

# Conclusions

If you have a Bayesian problem, then just

Do It !