

# Fast Quantification of Uncertainty and Robustness with Variational Bayes

Tamara Broderick  
ITT Career Development  
Assistant Professor,  
MIT

With: Ryan Giordano, Rachael Meager, Jonathan H. Huggins, Michael I. Jordan

- Bayesian inference

- Bayesian inference
  - Complex, modular models

- Bayesian inference
  - Complex, modular models; posterior distribution

- Bayesian inference  $p(\theta)$ 
  - Complex, modular models; posterior distribution

- Bayesian inference  $p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution

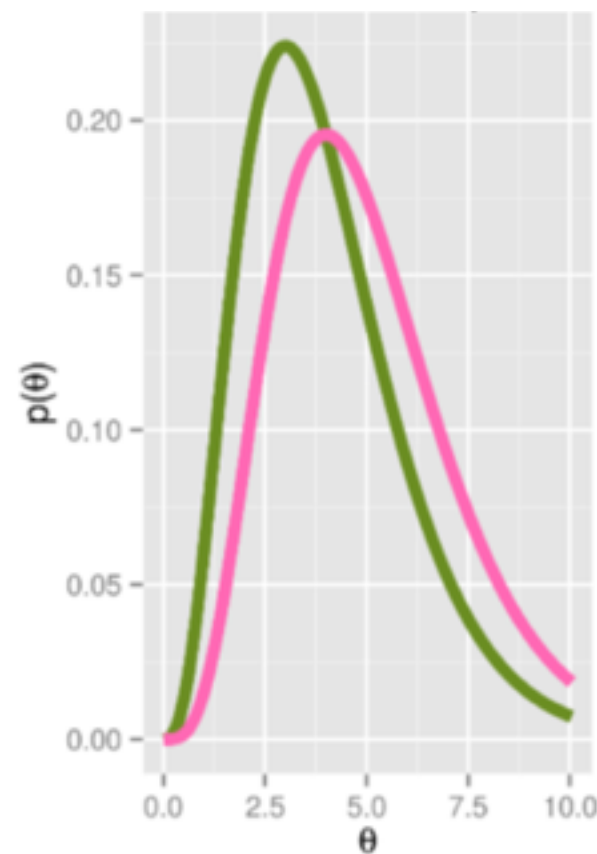


- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective

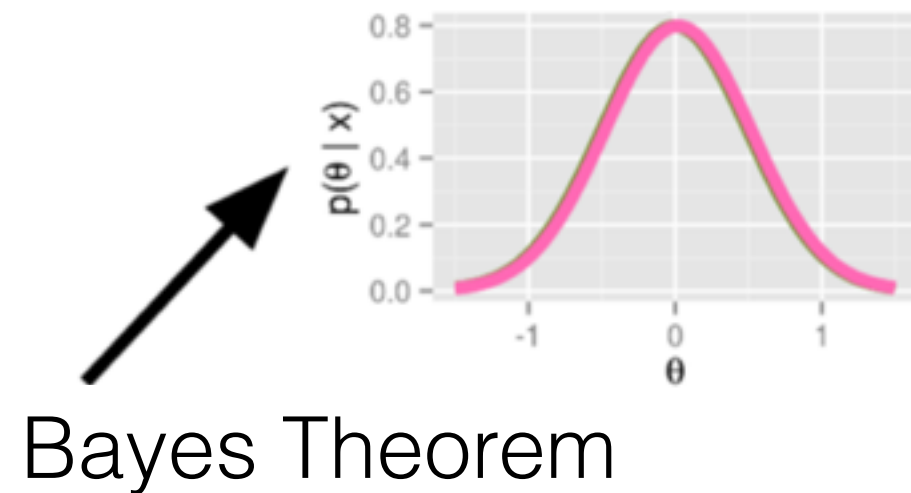
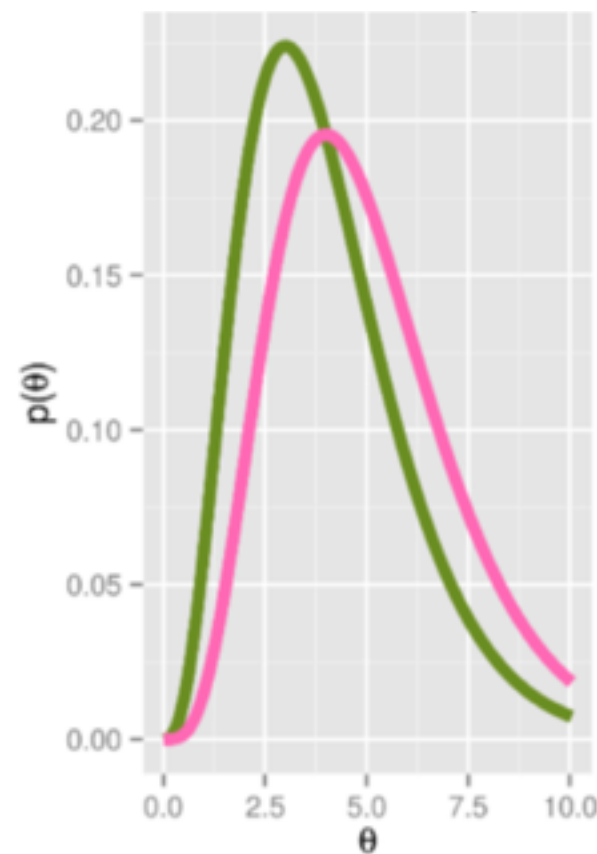
- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective

Some reasonable priors



- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective

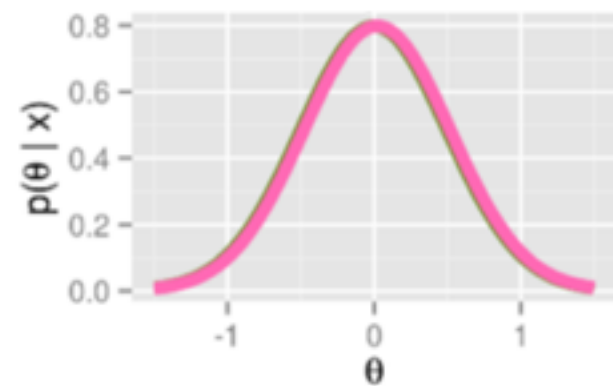
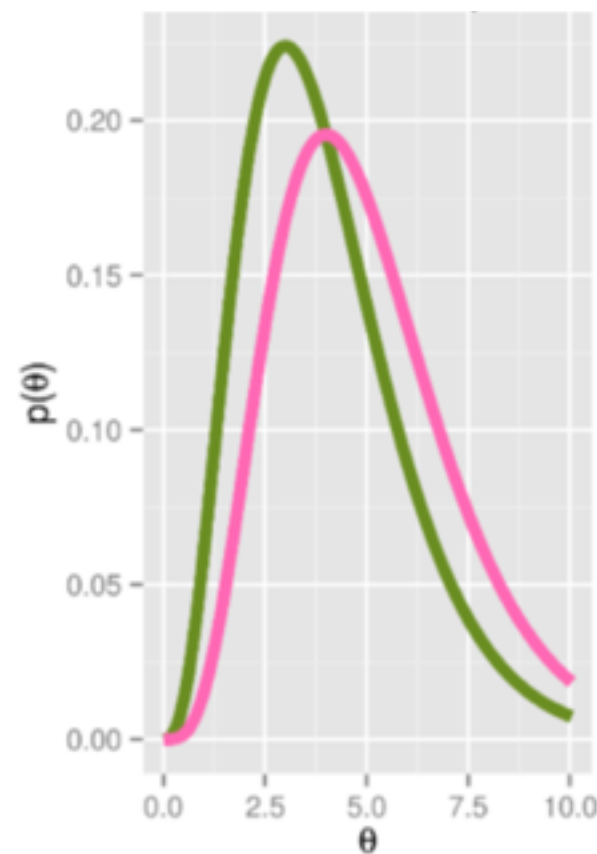
Some reasonable priors



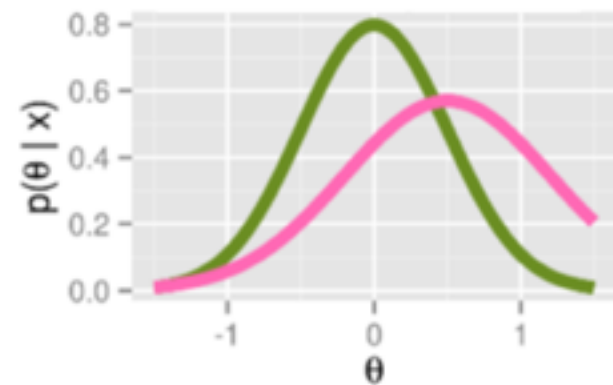
Bayes Theorem

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective

Some reasonable priors

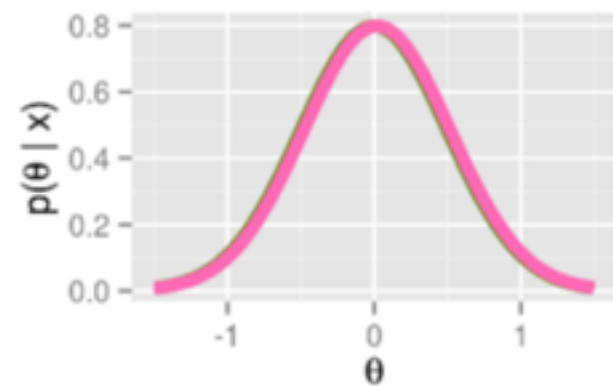
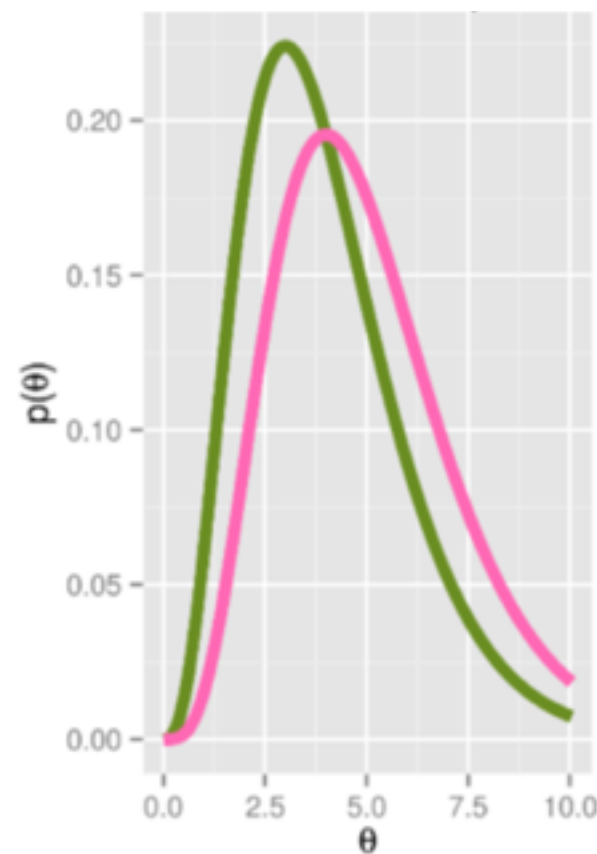


Bayes Theorem

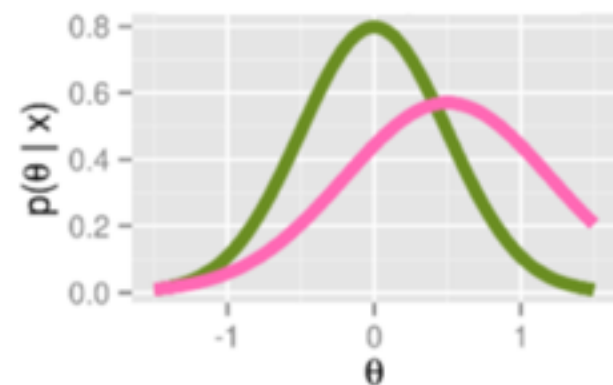


- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective; complex models

Some reasonable priors



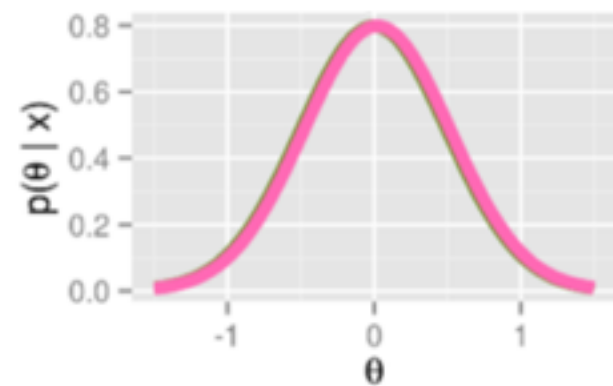
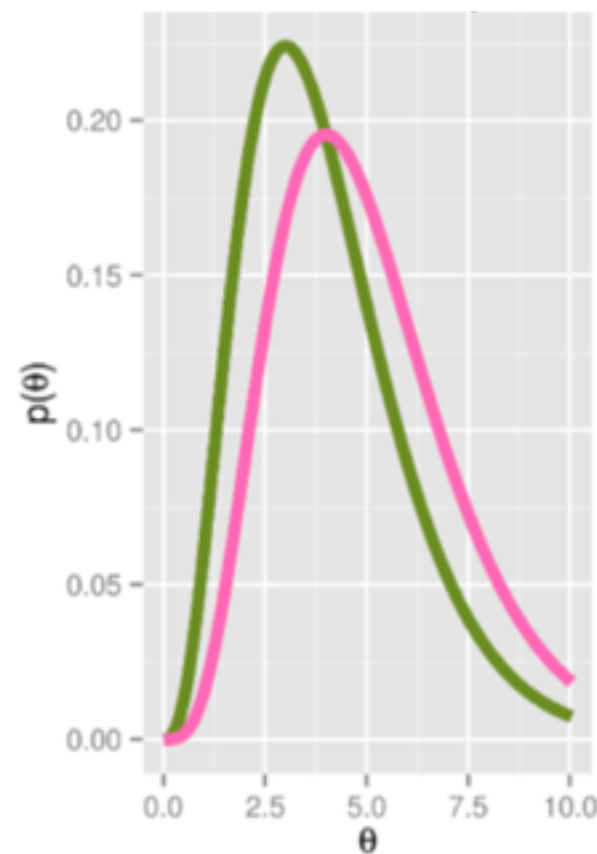
Bayes Theorem



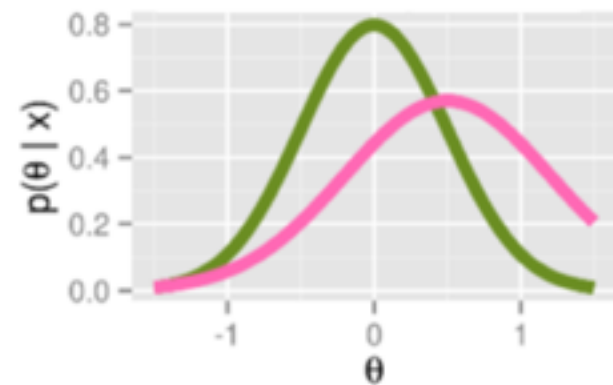
# robustness quantification

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective; complex models

Some reasonable priors



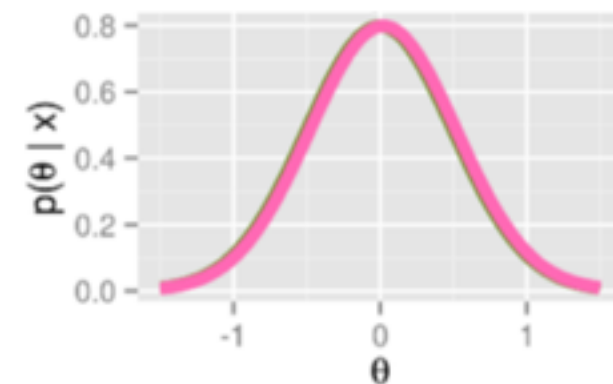
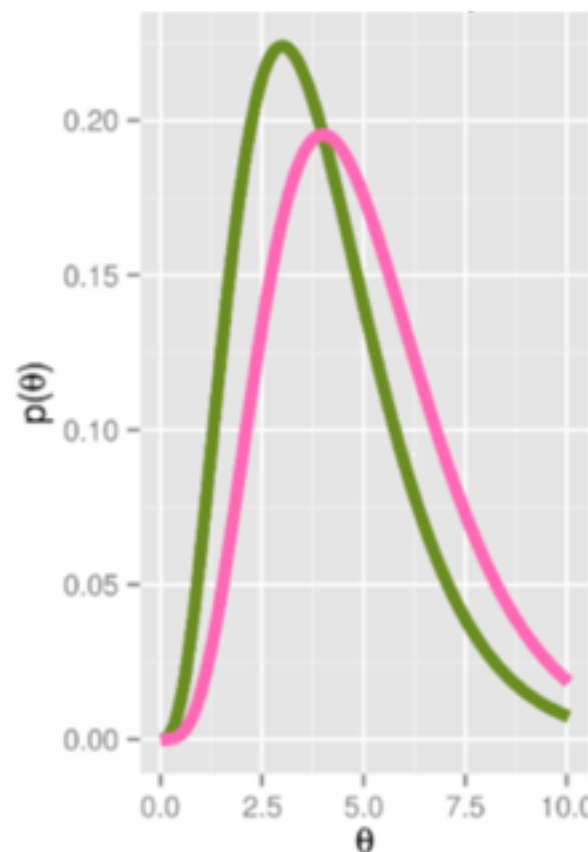
Bayes Theorem



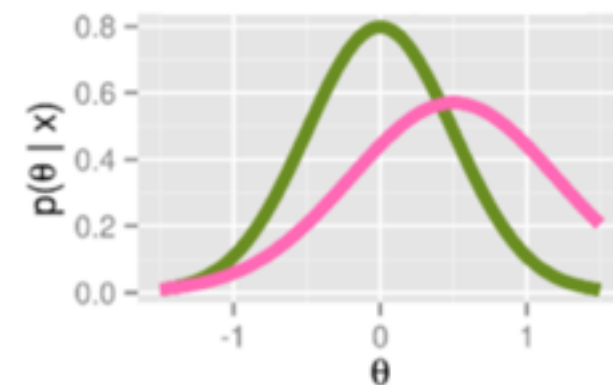
# robustness quantification

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective; complex models
- Challenge: Approximating the posterior can be computationally expensive

Some reasonable priors



Bayes Theorem

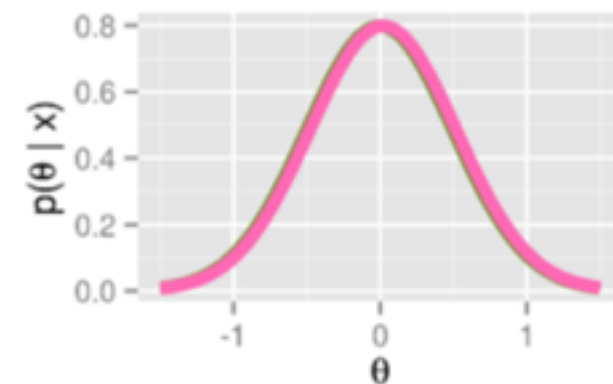
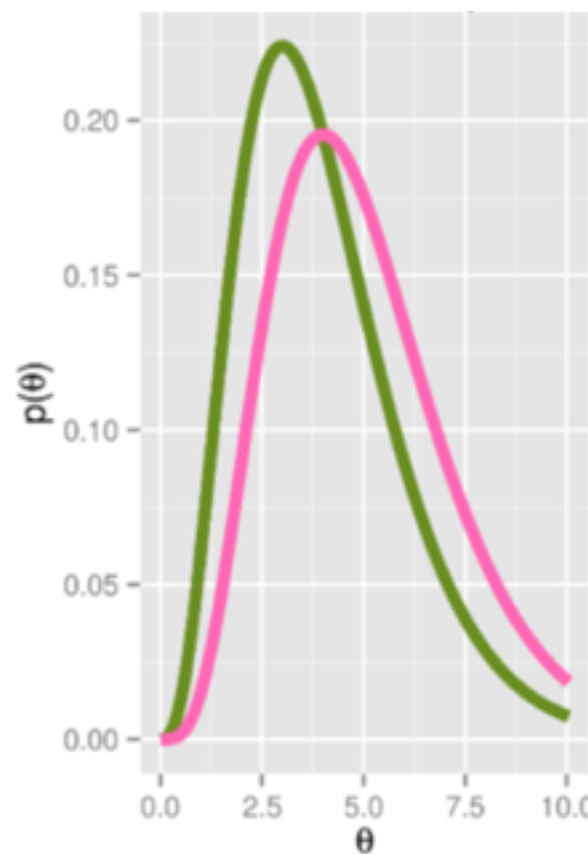




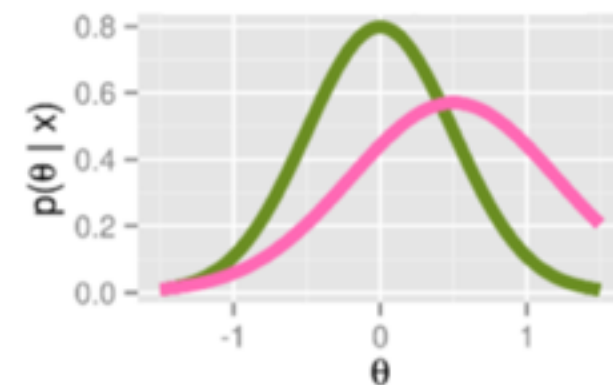
# robustness quantification

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective; complex models
- Challenge: Approximating the posterior can be computationally expensive

Some reasonable priors



Bayes Theorem

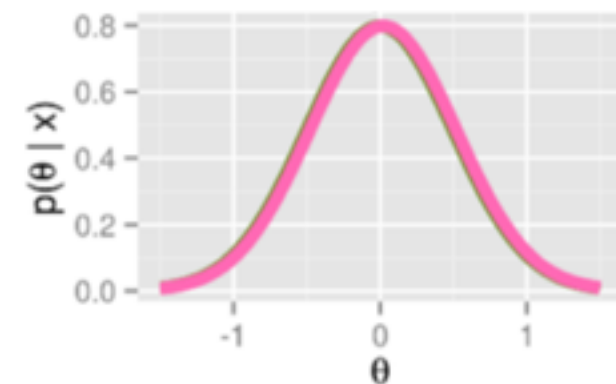
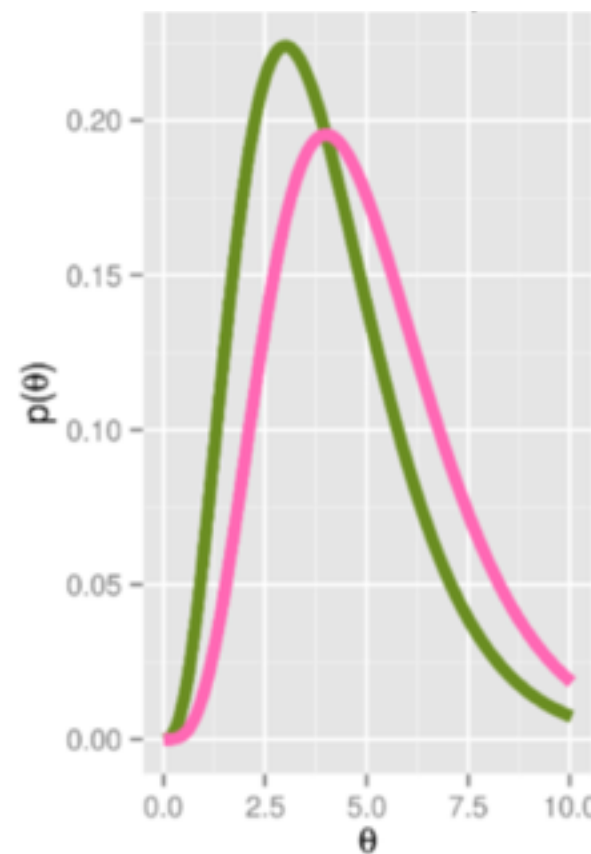


*variational Bayes*

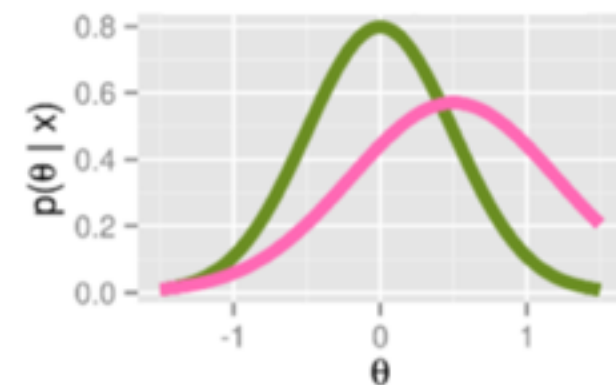
# Uncertainty & robustness quantification

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective; complex models
- Challenge: Approximating the posterior can be computationally expensive

Some reasonable priors



Bayes Theorem

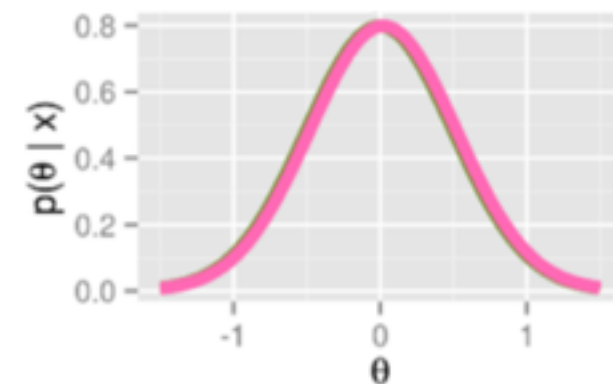
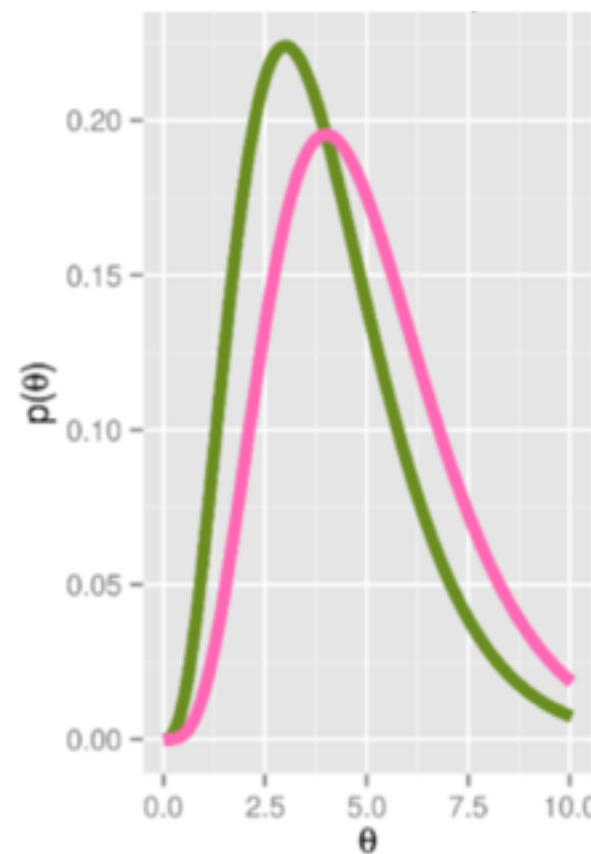


*variational Bayes*

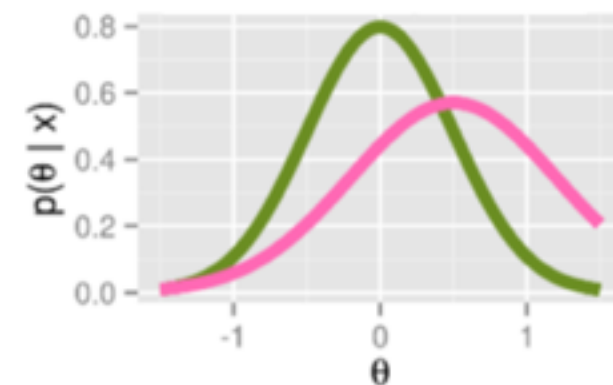
# Uncertainty & robustness quantification

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective; complex models
- Challenge: Approximating the posterior can be computationally expensive

Some reasonable priors



Bayes Theorem

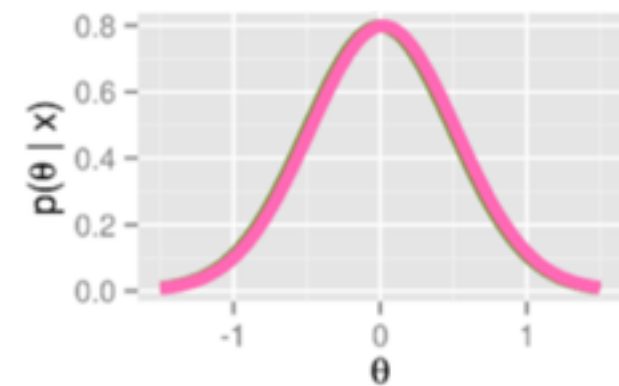
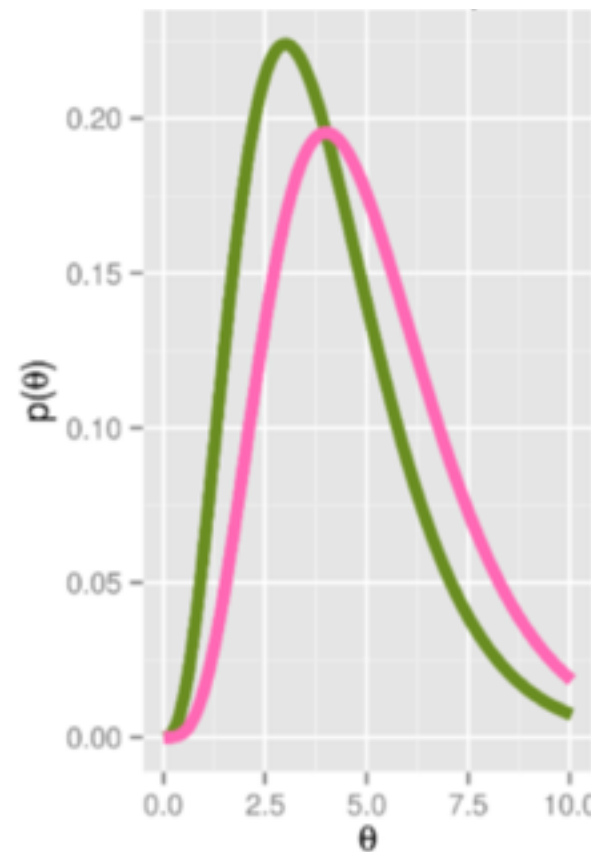


- We propose: *linear response variational Bayes*

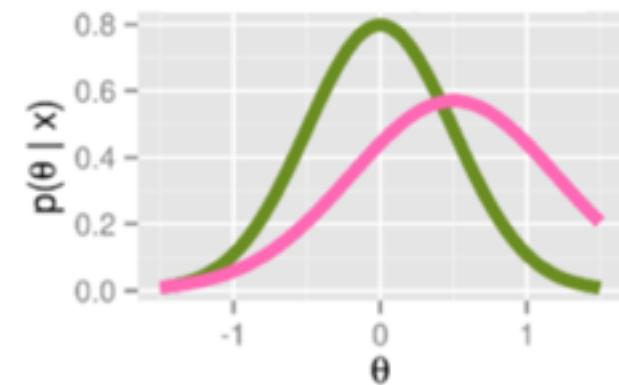
# Uncertainty & robustness quantification

- Bayesian inference  $p(\theta|x) \propto_{\theta} p(x|\theta)p(\theta)$ 
  - Complex, modular models; posterior distribution
- Challenge: Express prior beliefs in a distribution
  - Time-consuming; subjective; complex models
- Challenge: Approximating the posterior can be computationally expensive

Some reasonable priors



Bayes Theorem



- We propose: *linear response variational Bayes*

[see also Oppen, Winther 2003]

# Roadmap

# Roadmap

- Variational Bayes as an alternative to MCMC

# Roadmap

- Variational Bayes as an alternative to MCMC
- Challenges of VB

# Roadmap

- Variational Bayes as an alternative to MCMC
- Challenges of VB
- Accurate uncertainties from VB



# Roadmap

- Variational Bayes as an alternative to MCMC
- Challenges of VB
- Accurate uncertainties from VB
- Accurate robustness quantification from VB

# Roadmap

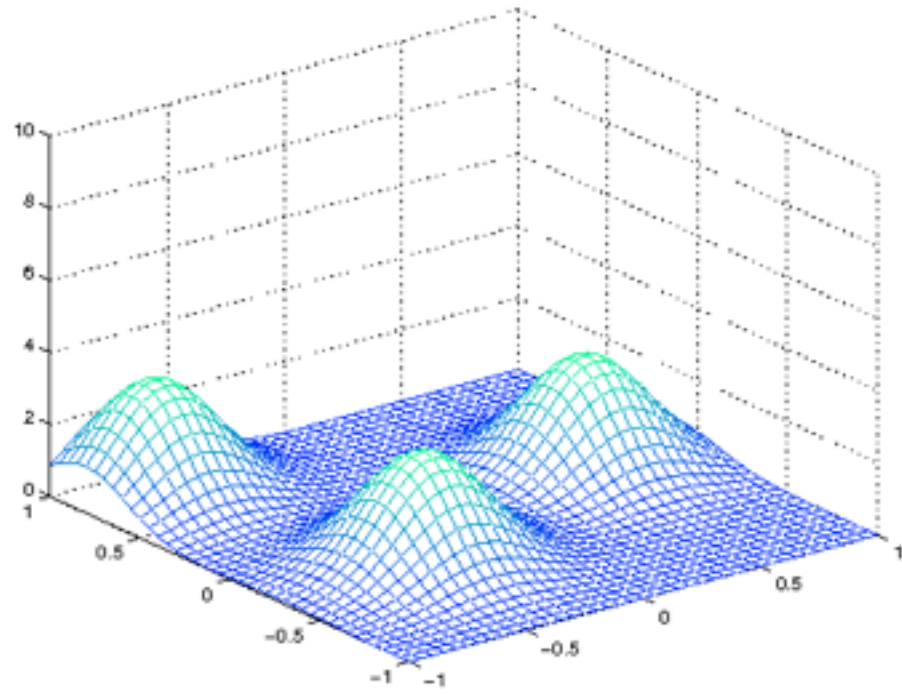
- Variational Bayes as an alternative to MCMC
  - Challenges of VB
  - Accurate uncertainties from VB
  - Accurate robustness quantification from VB
- 
- Big idea: derivatives/perturbations are relatively easy in VB

# Variational Bayes

# Variational Bayes

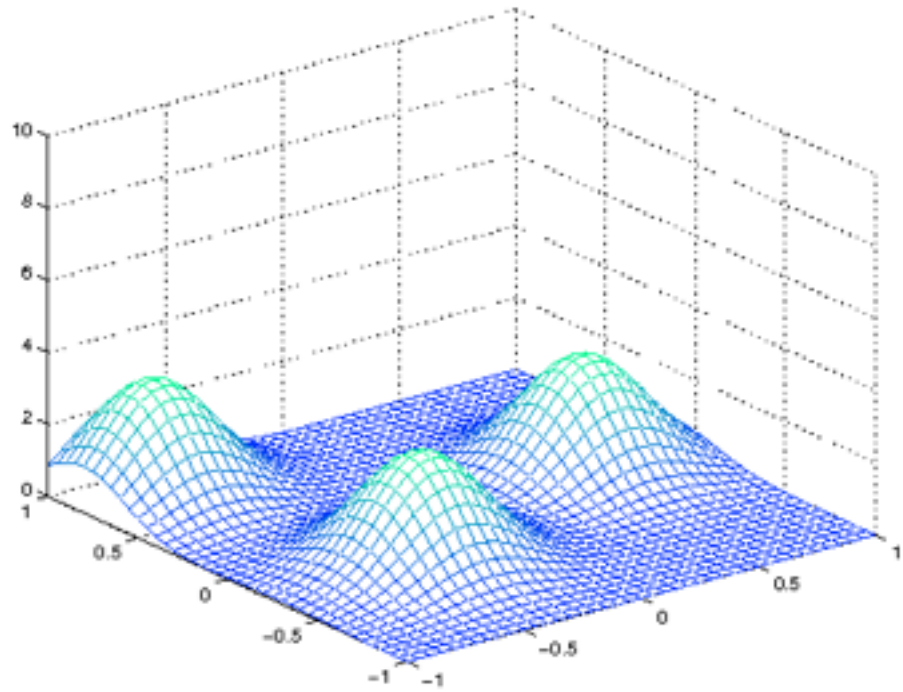
- VB approximation

# Variational Bayes



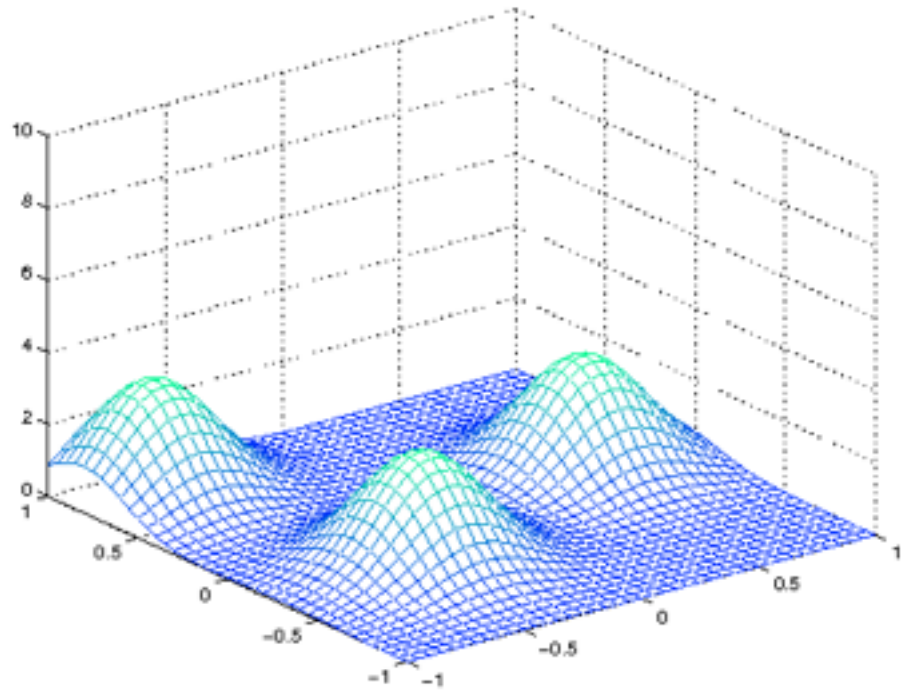
- VB approximation

# Variational Bayes

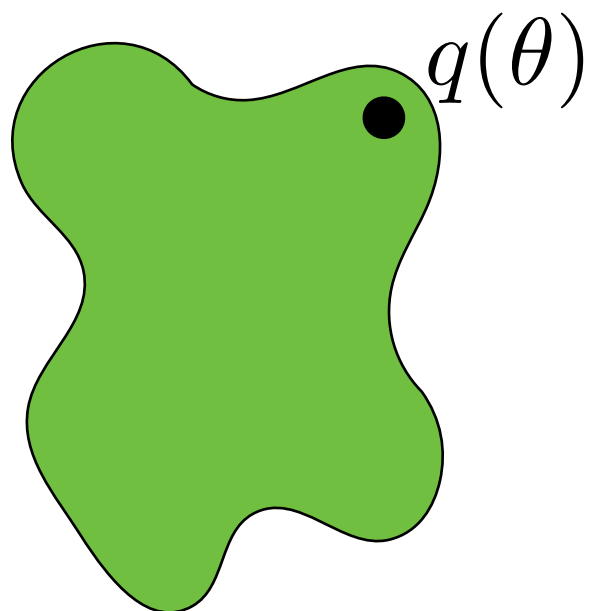


- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$

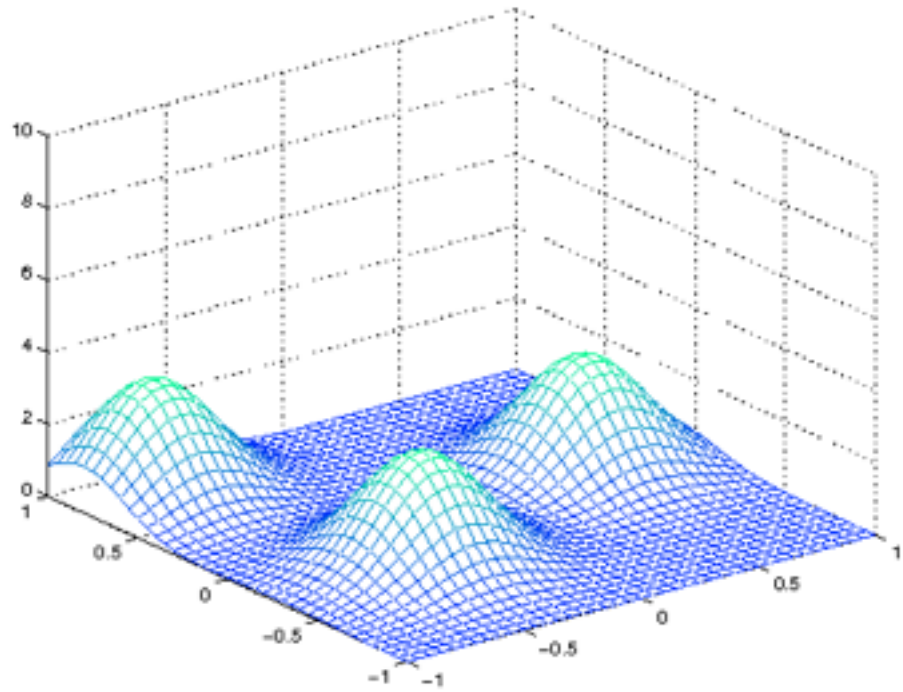
# Variational Bayes



- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$

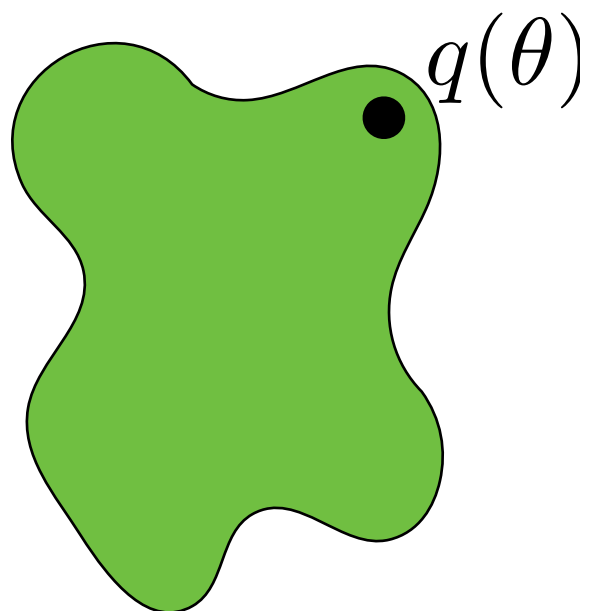


# Variational Bayes



- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$

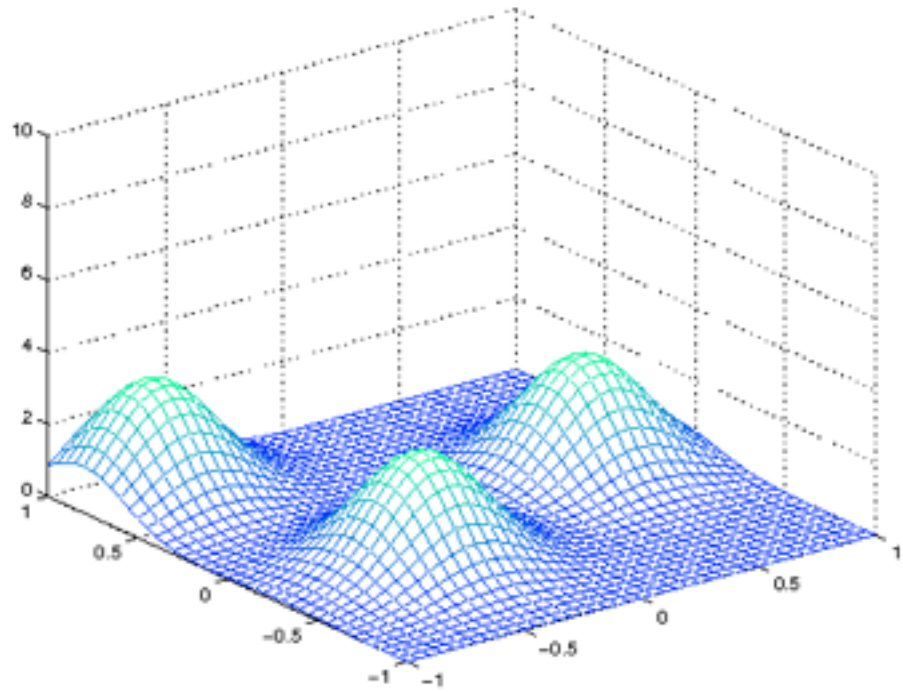
$p(\theta|x)$



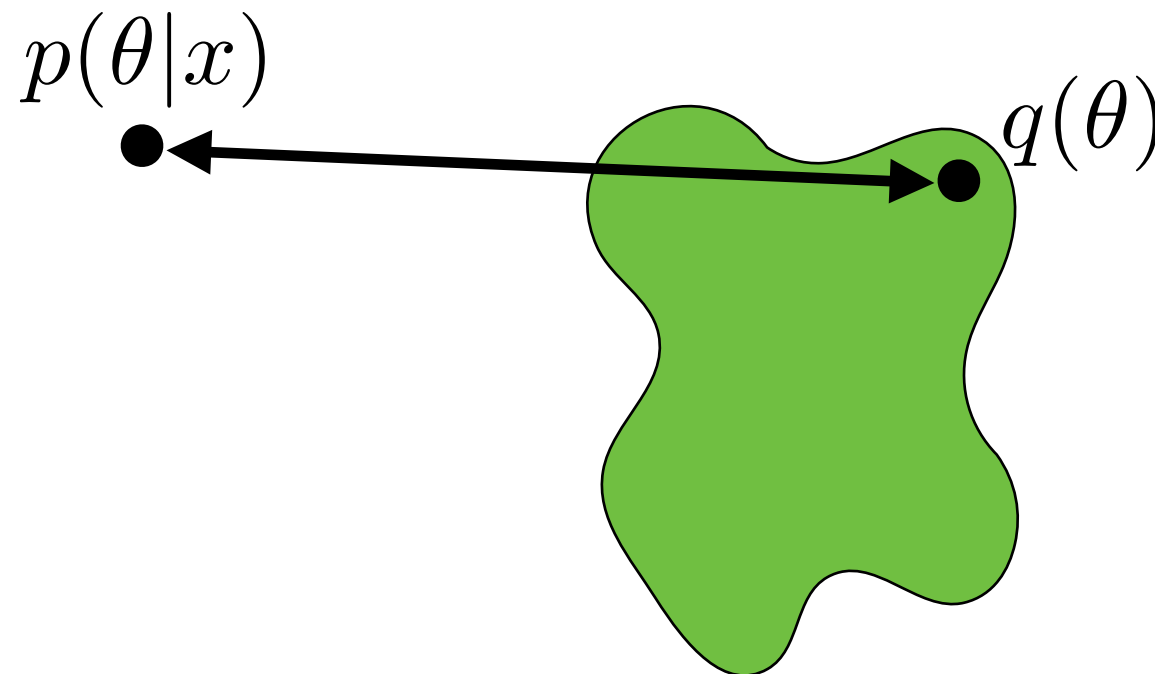
$q(\theta)$



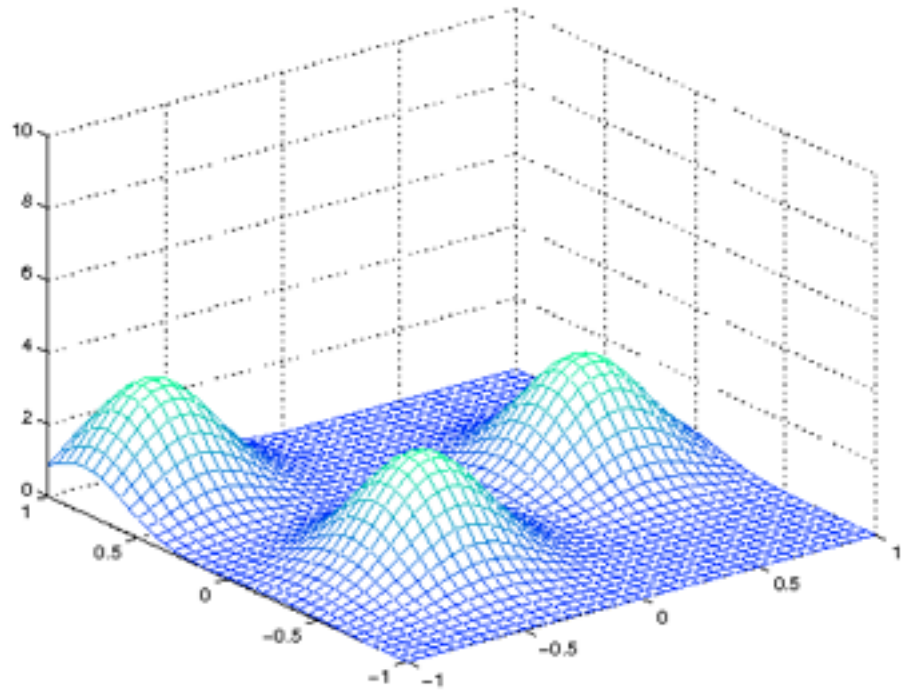
# Variational Bayes



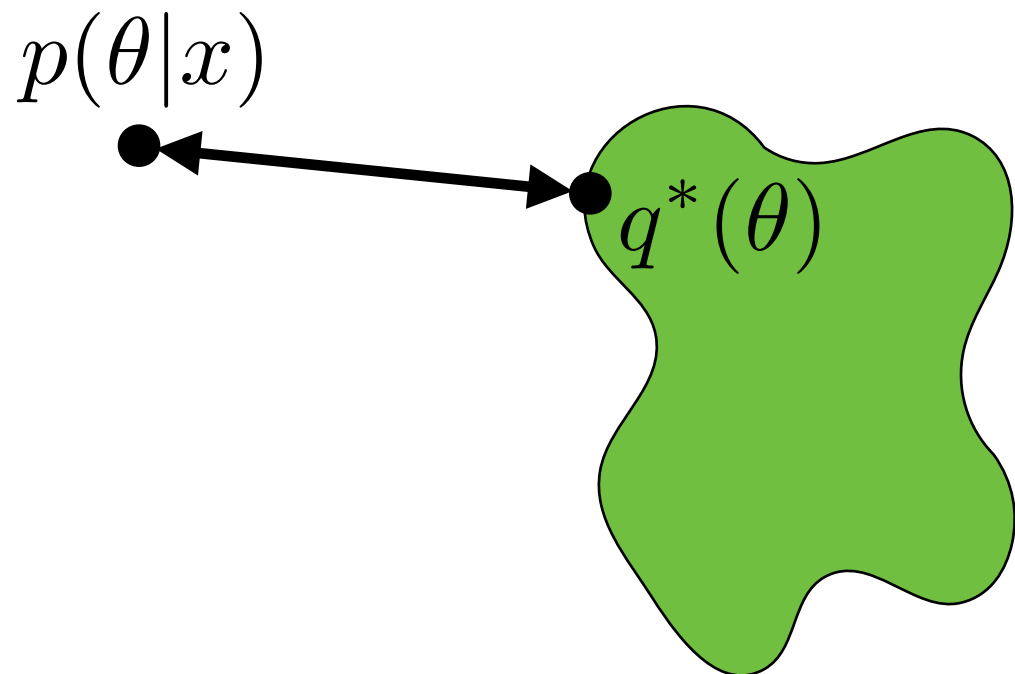
- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$



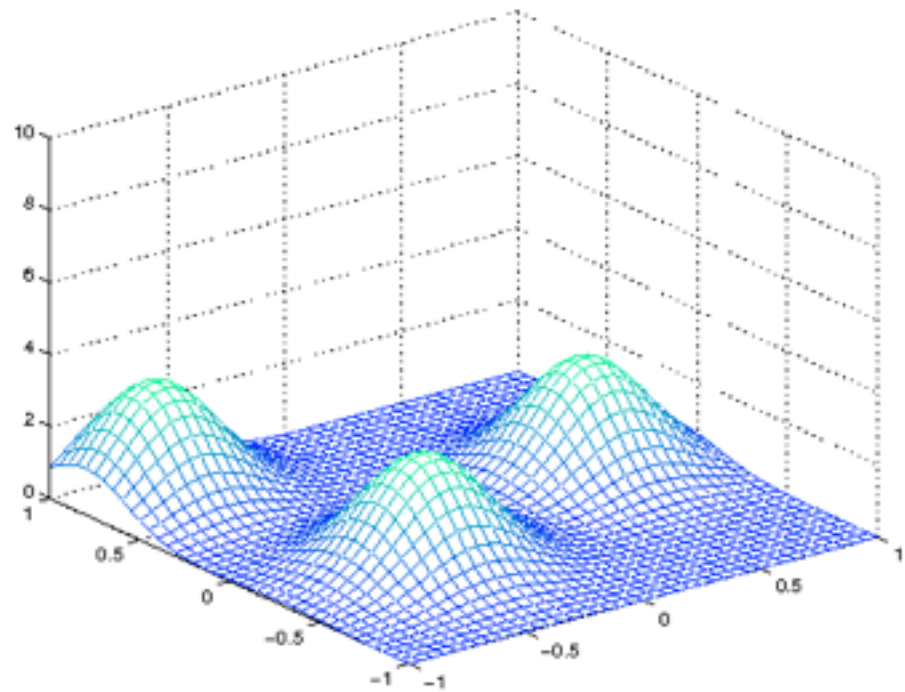
# Variational Bayes



- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$

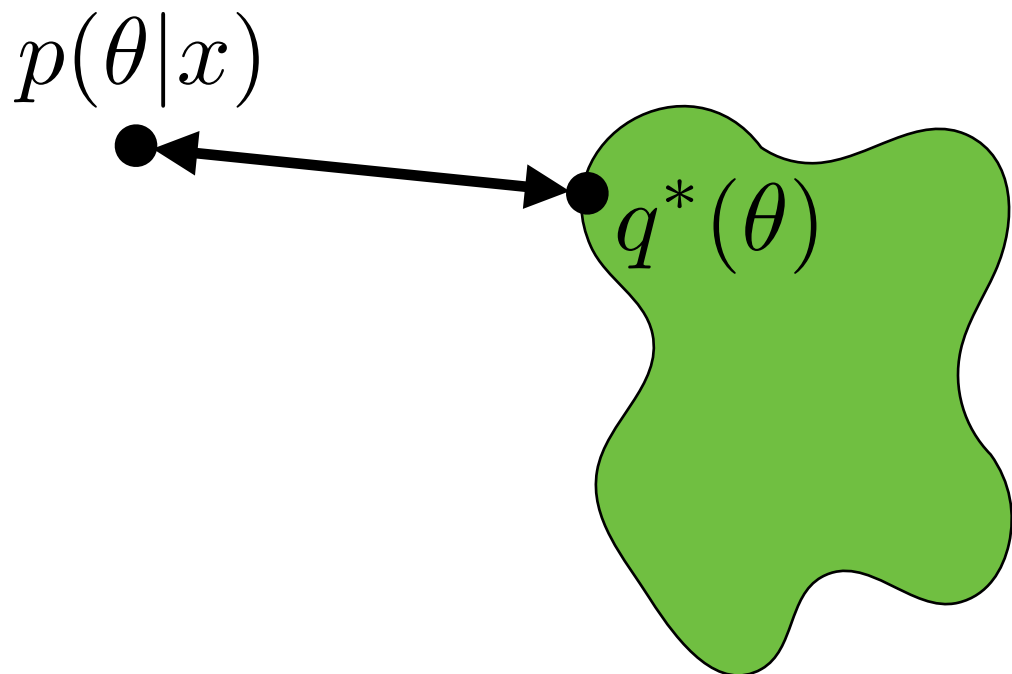


# Variational Bayes

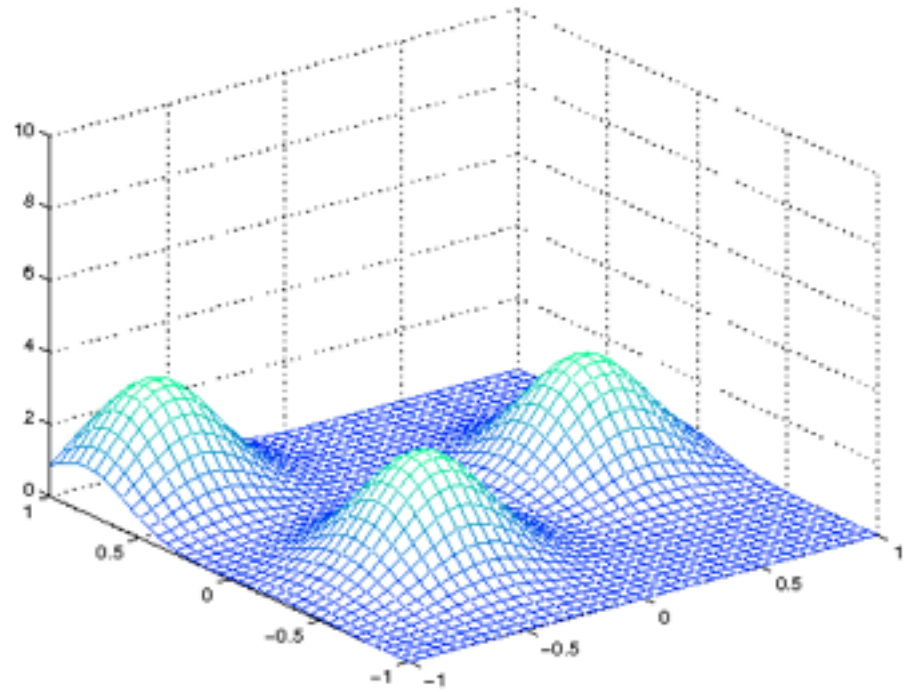


- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - Minimize Kullback-Leibler (KL) divergence:

$$KL(q||p(\cdot|x))$$

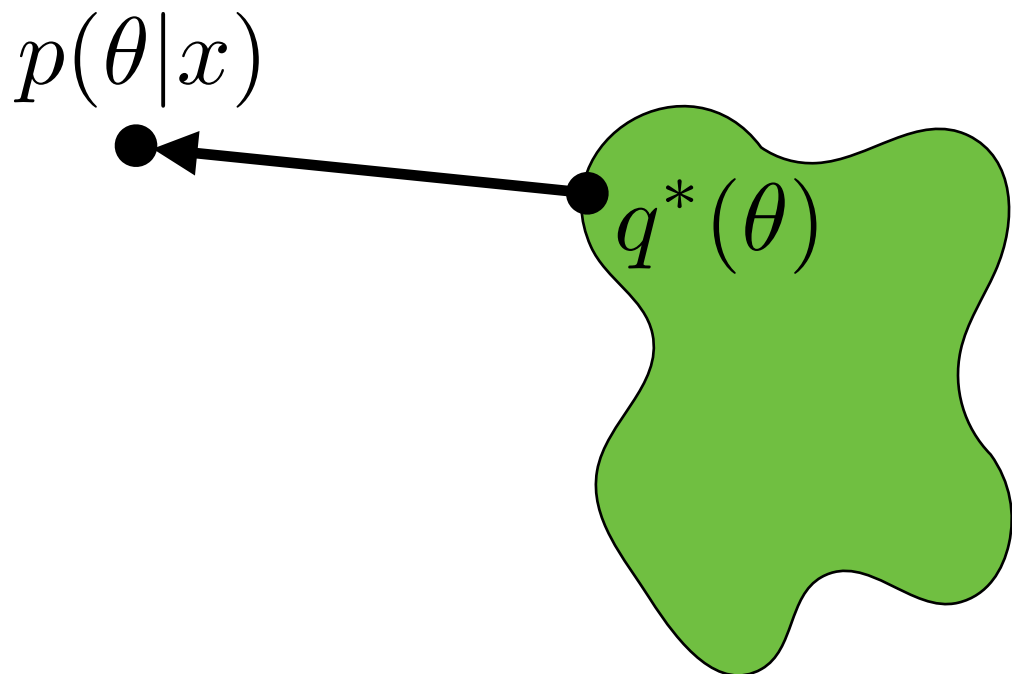


# Variational Bayes

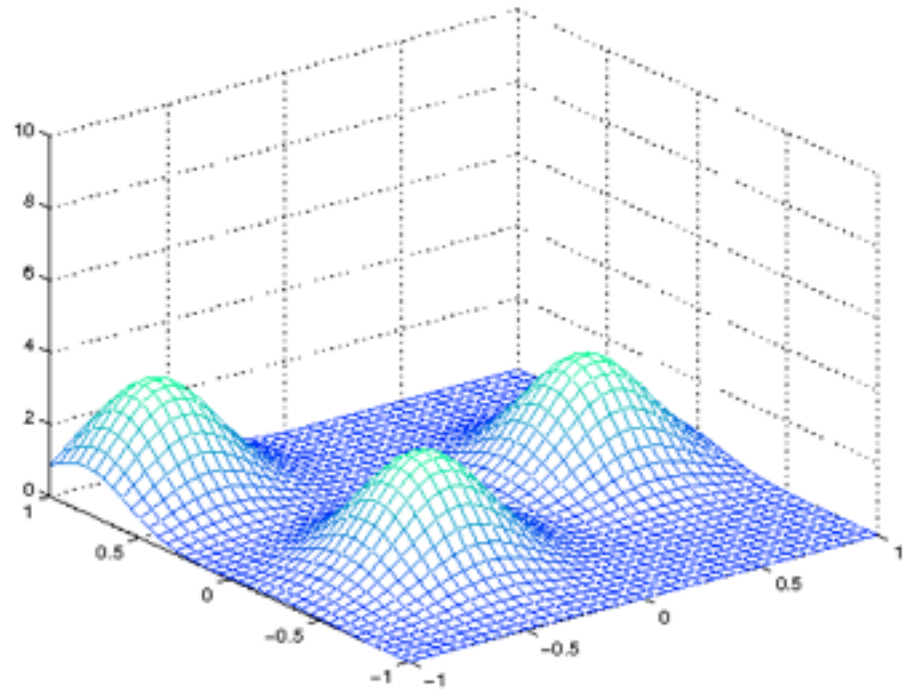


- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - Minimize Kullback-Leibler (KL) divergence:

$$KL(q||p(\cdot|x))$$

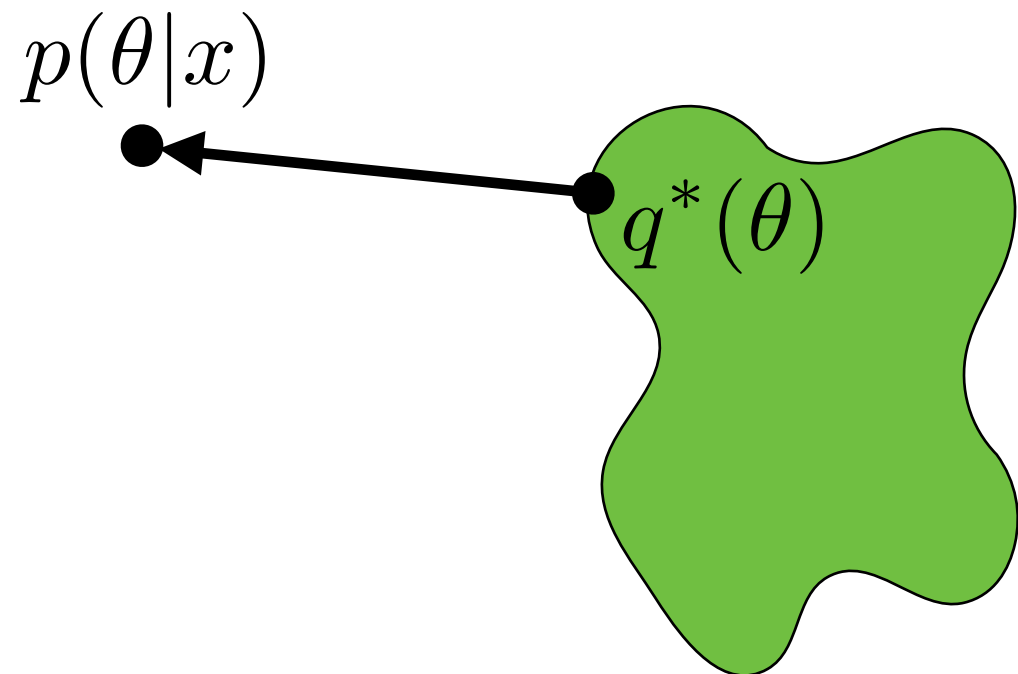


# Variational Bayes



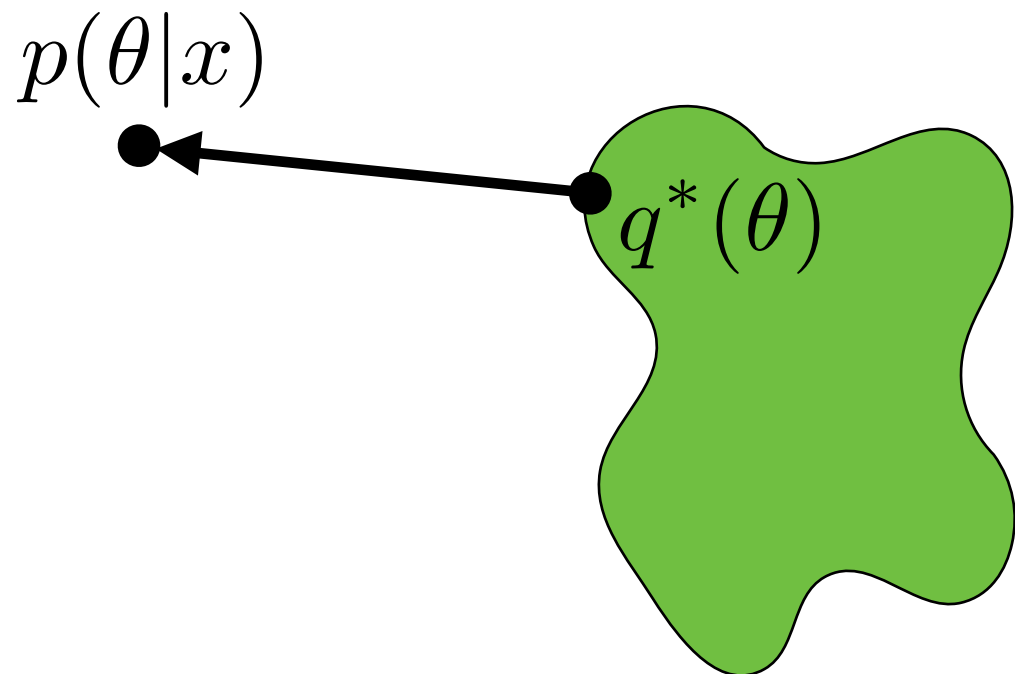
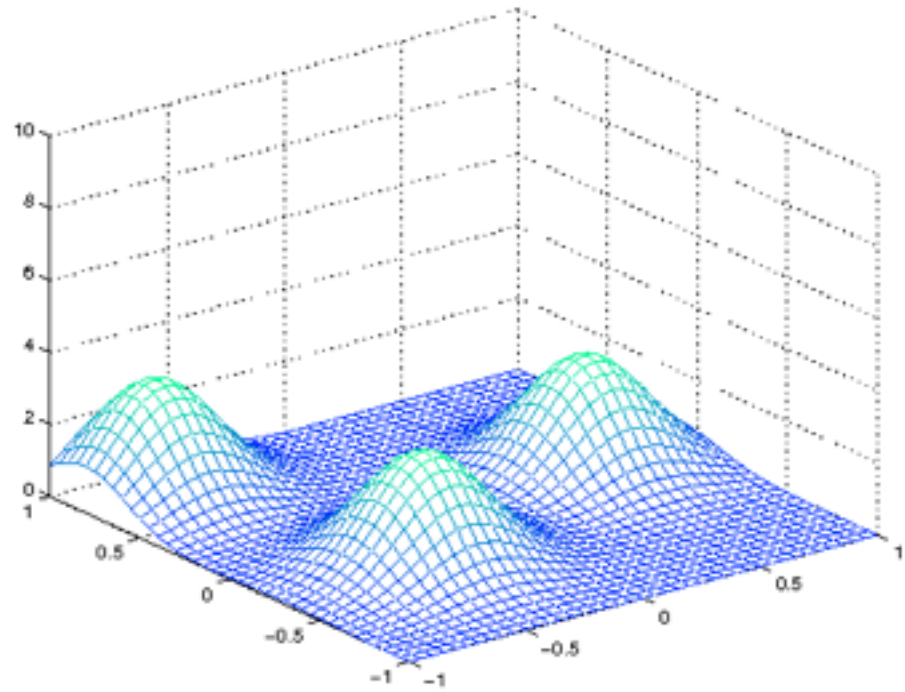
- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - Minimize Kullback-Leibler (KL) divergence:

$$KL(q||p(\cdot|x))$$



- VB practical success

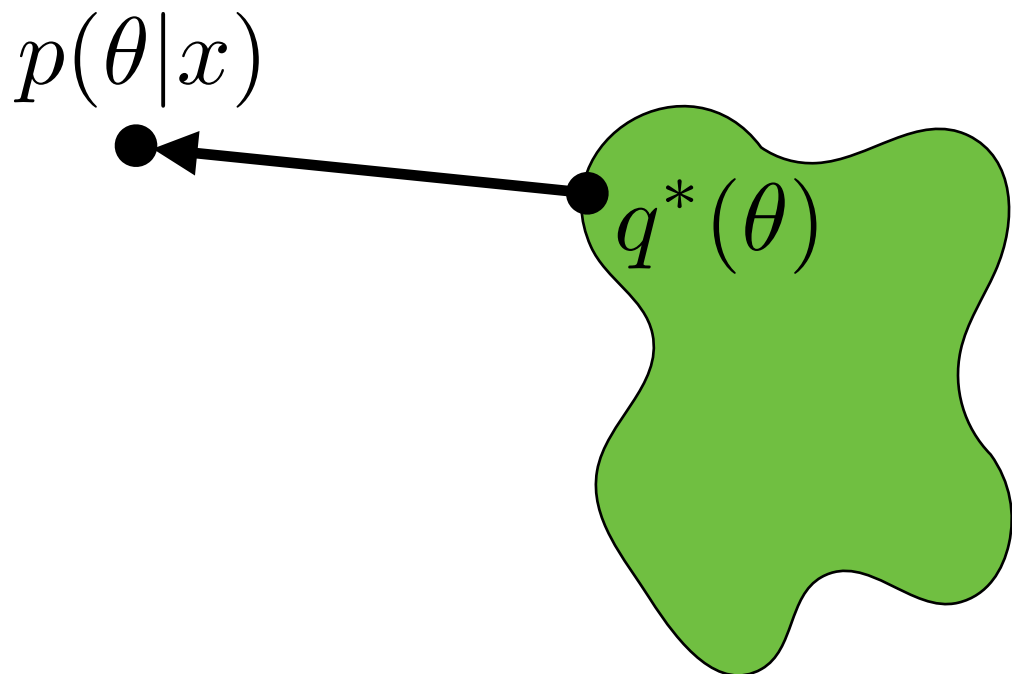
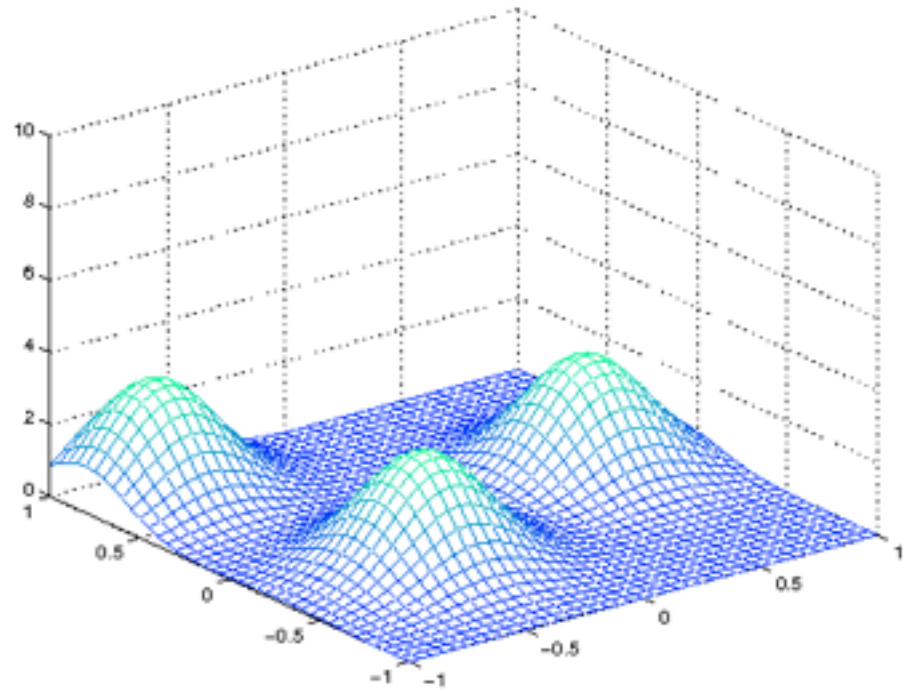
# Variational Bayes



- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - Minimize Kullback-Leibler (KL) divergence:
$$KL(q||p(\cdot|x))$$
- VB practical success
  - point estimates and prediction

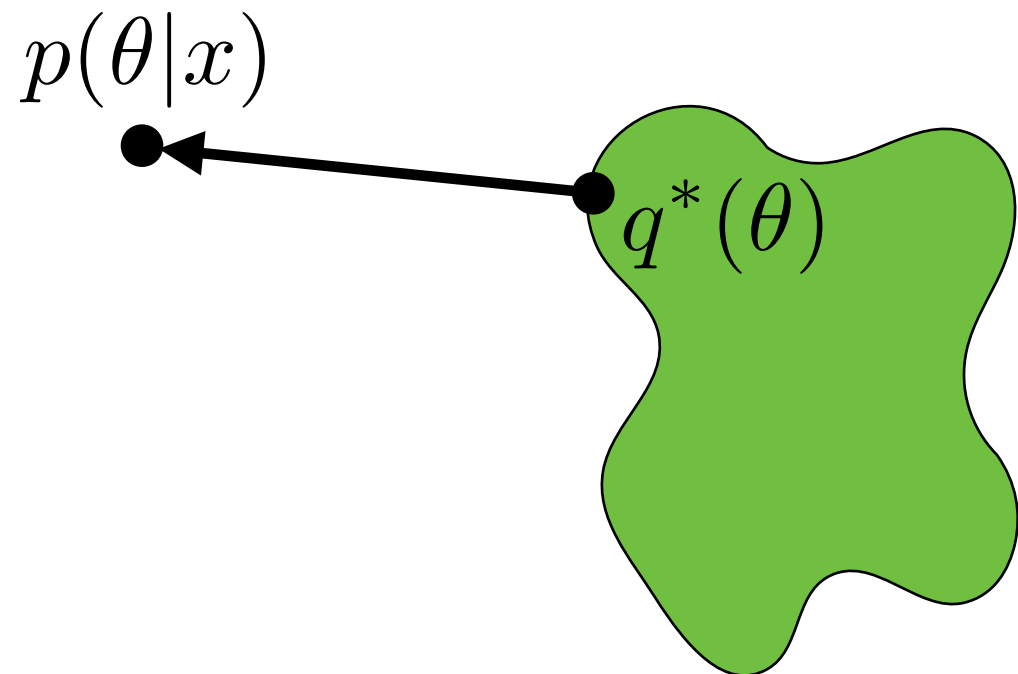
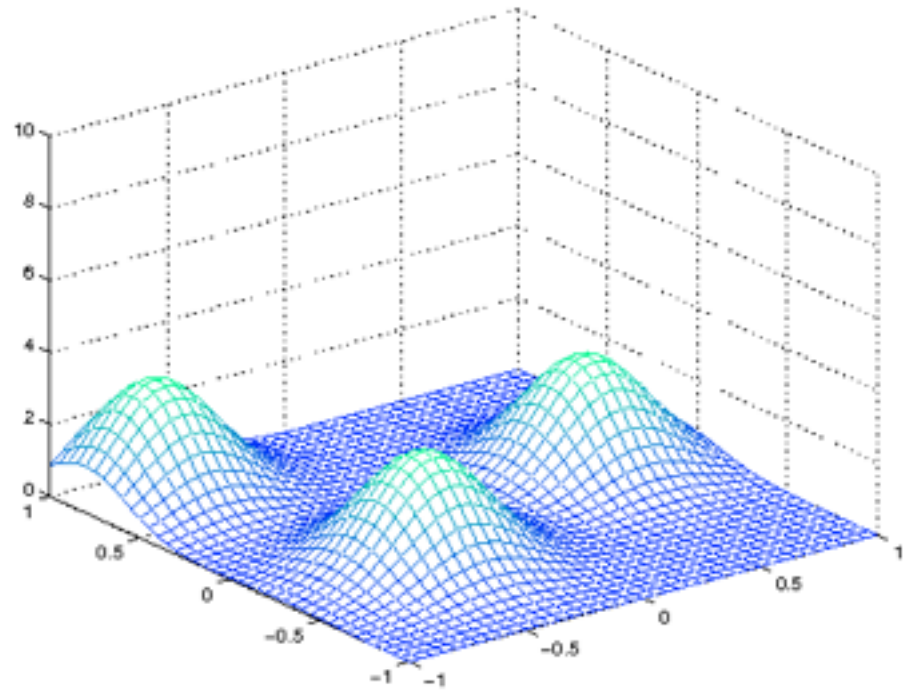


# Variational Bayes



- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - Minimize Kullback-Leibler (KL) divergence:
$$KL(q||p(\cdot|x))$$
- VB practical success
  - point estimates and prediction
  - fast

# Variational Bayes



- VB approximation
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - Minimize Kullback-Leibler (KL) divergence:
$$KL(q||p(\cdot|x))$$
- VB practical success
  - point estimates and prediction
  - fast, streaming, distributed



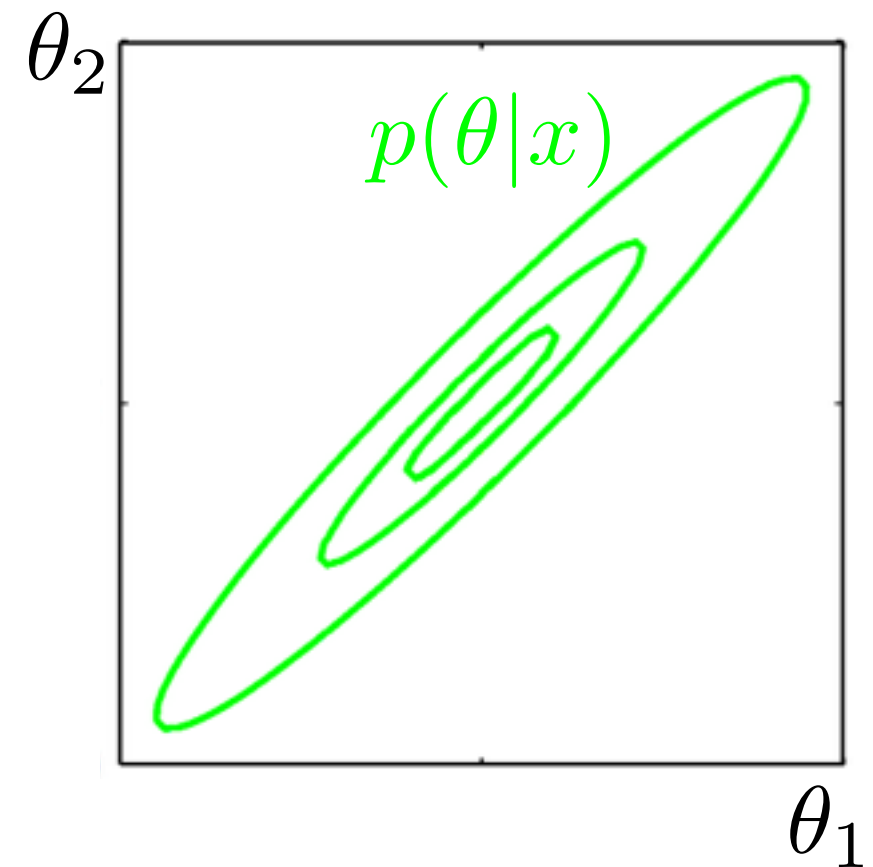
# What about uncertainty?

# What about uncertainty?

- Variational Bayes

# What about uncertainty?

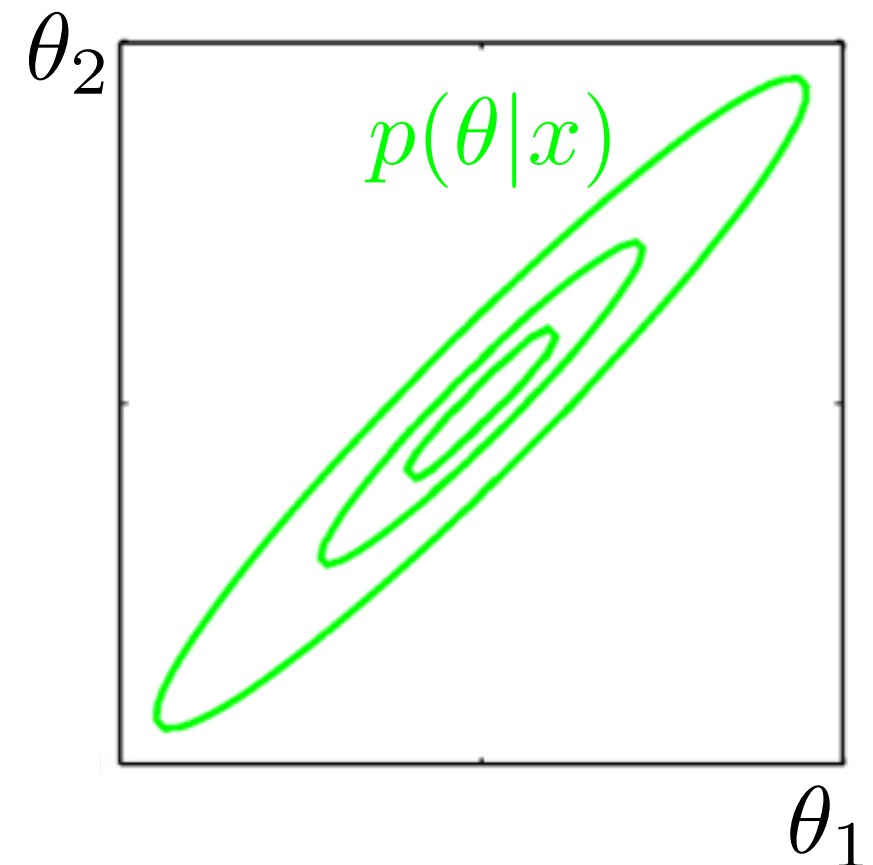
- Variational Bayes



# What about uncertainty?

- Variational Bayes

$$KL(q||p(\cdot|x)) = \int_{\theta} q(\theta) \log \frac{q(\theta)}{p(\theta|x)} d\theta$$



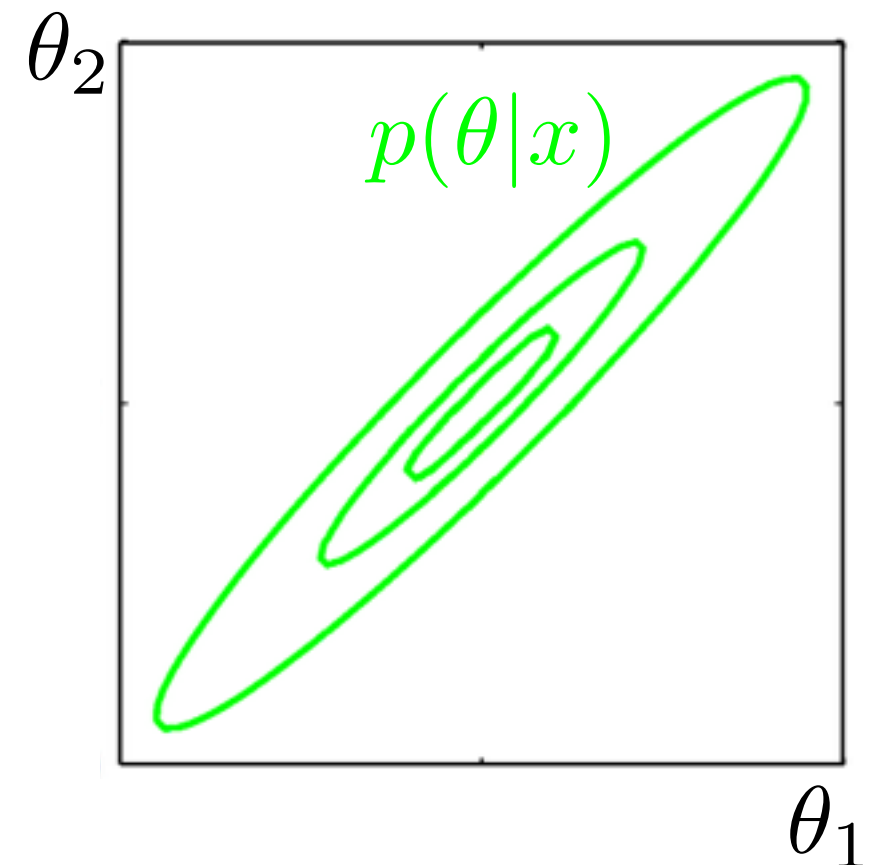
# What about uncertainty?

- Variational Bayes

$$KL(q||p(\cdot|x)) = \int_{\theta} q(\theta) \log \frac{q(\theta)}{p(\theta|x)} d\theta$$

- Mean-field variational Bayes (MFVB)

$$q(\theta) = \prod_{j=1}^J q(\theta_j)$$



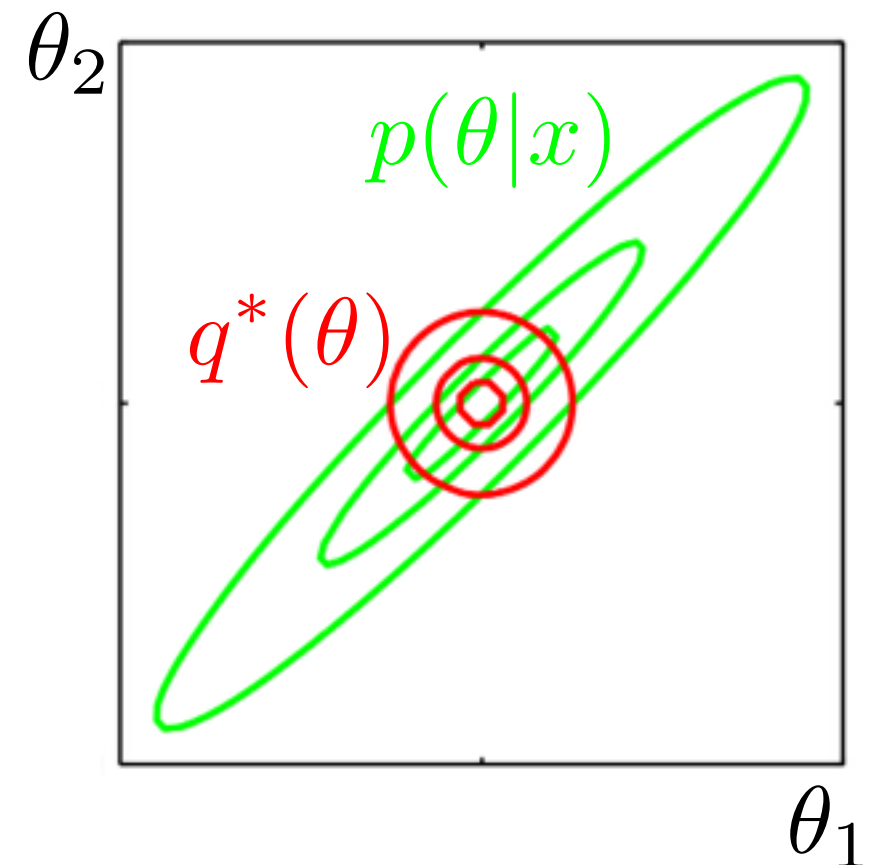
# What about uncertainty?

- Variational Bayes

$$KL(q||p(\cdot|x)) = \int_{\theta} q(\theta) \log \frac{q(\theta)}{p(\theta|x)} d\theta$$

- Mean-field variational Bayes (MFVB)

$$q(\theta) = \prod_{j=1}^J q(\theta_j)$$



# What about uncertainty?

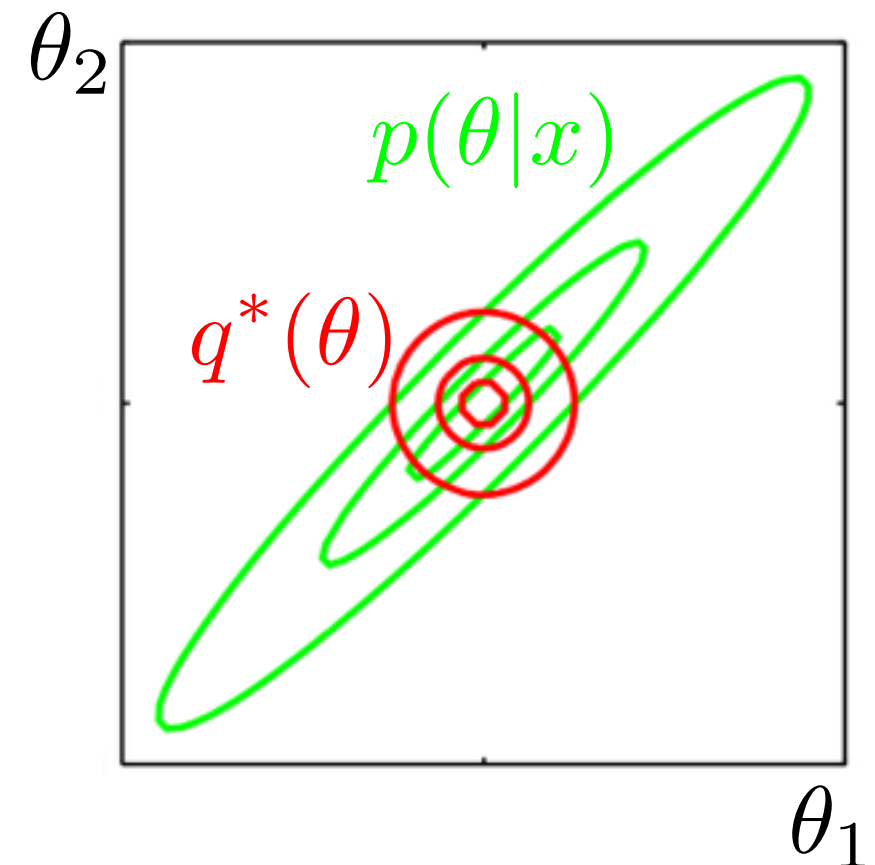
- Variational Bayes

$$KL(q||p(\cdot|x)) = \int_{\theta} q(\theta) \log \frac{q(\theta)}{p(\theta|x)} d\theta$$

- Mean-field variational Bayes (MFVB)

$$q(\theta) = \prod_{j=1}^J q(\theta_j)$$

- Underestimates variance (sometimes severely)



# What about uncertainty?

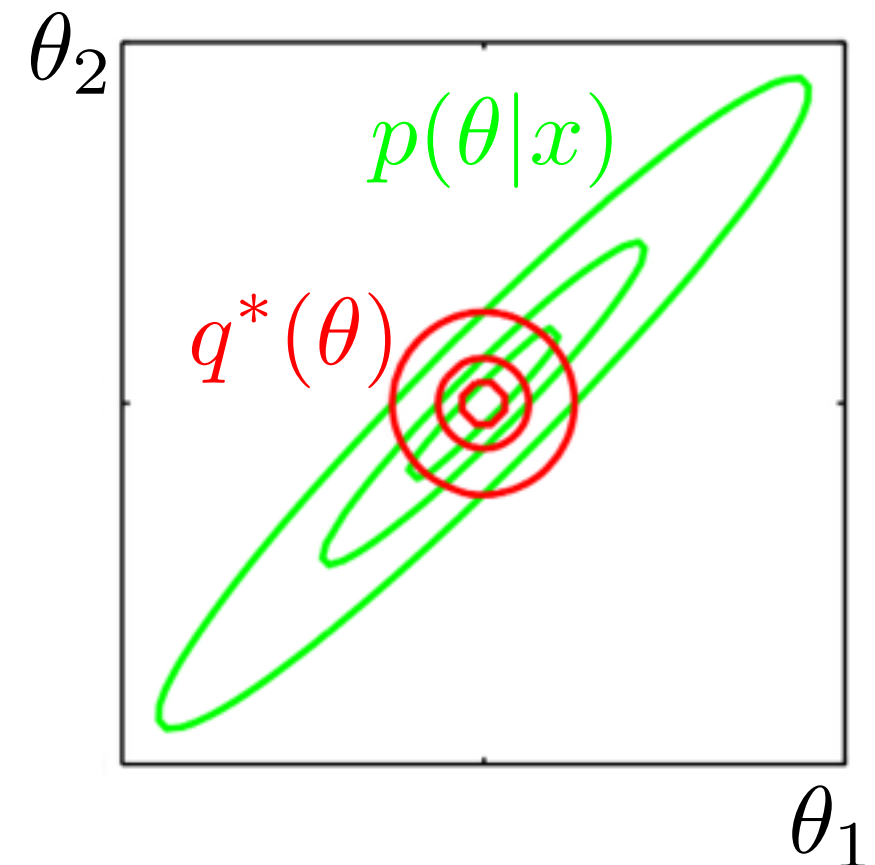
- Variational Bayes

$$KL(q||p(\cdot|x)) = \int_{\theta} q(\theta) \log \frac{q(\theta)}{p(\theta|x)} d\theta$$

- Mean-field variational Bayes (MFVB)

$$q(\theta) = \prod_{j=1}^J q(\theta_j)$$

- Underestimates variance (sometimes severely)
- No covariance estimates





# What about uncertainty?

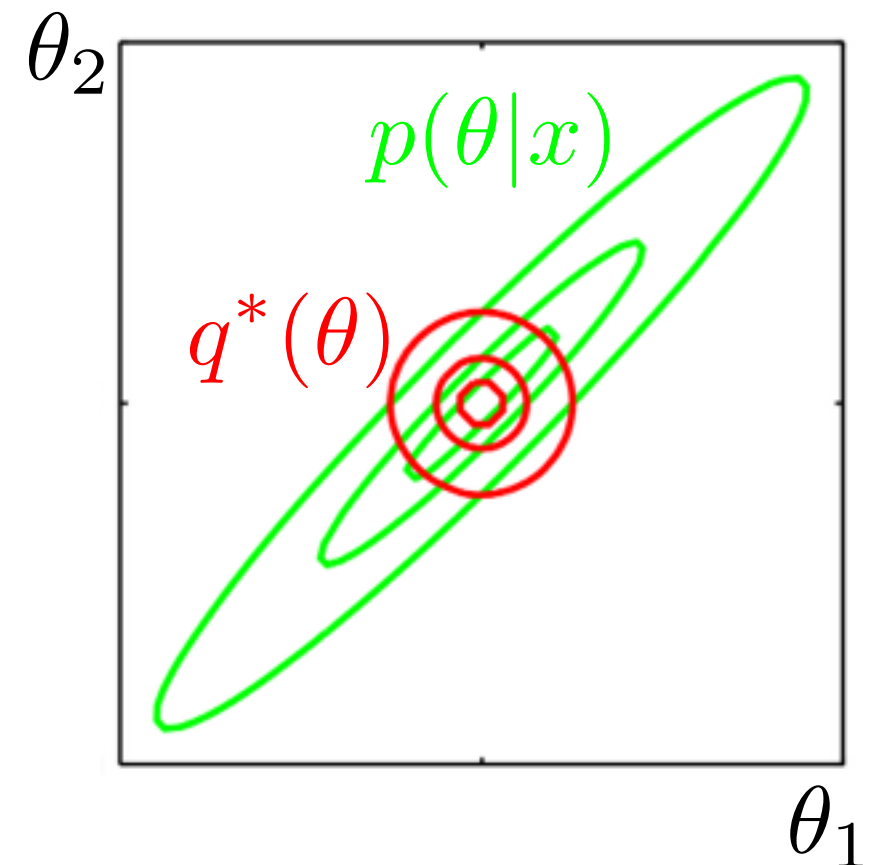
- Variational Bayes

$$KL(q||p(\cdot|x)) = \int_{\theta} q(\theta) \log \frac{q(\theta)}{p(\theta|x)} d\theta$$

- Mean-field variational Bayes (MFVB)

$$q(\theta) = \prod_{j=1}^J q(\theta_j)$$

- Underestimates variance (sometimes severely)
- No covariance estimates



# What about uncertainty?

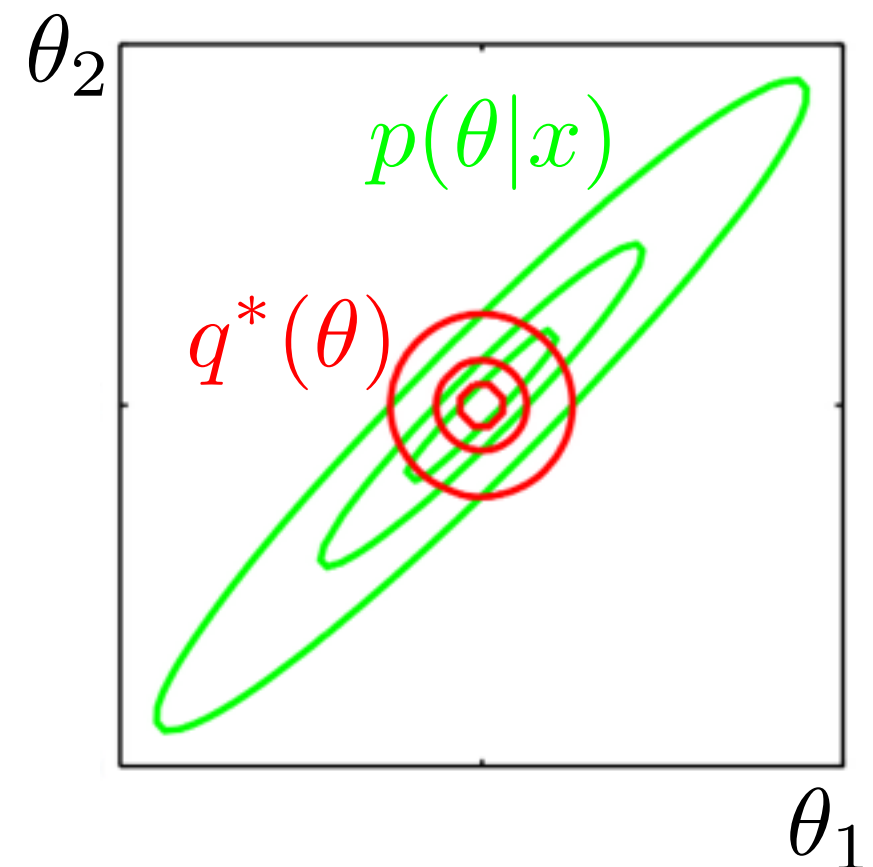
- Variational Bayes

$$KL(q||p(\cdot|x)) = \int_{\theta} q(\theta) \log \frac{q(\theta)}{p(\theta|x)} d\theta$$

- Mean-field variational Bayes (MFVB)

$$q(\theta) = \prod_{j=1}^J q(\theta_j)$$

- Underestimates variance (sometimes severely)
- No covariance estimates



[MacKay 2003; Bishop 2006; Wang, Titterton 2004; Turner, Sahani 2011]

[Fosdick 2013; Dunson 2014; Bardenet, Doucet, Holmes 2015]

# Linear response

# Linear response

- Cumulant-generating function

# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

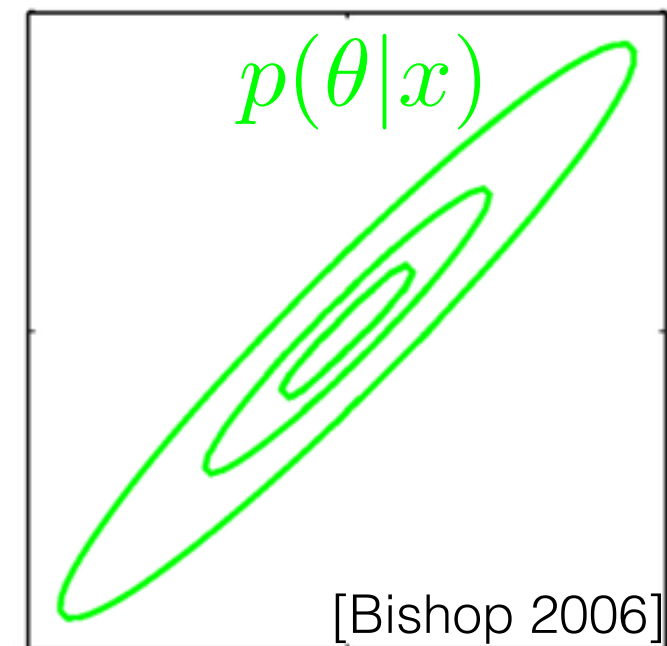
# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance



# Linear response

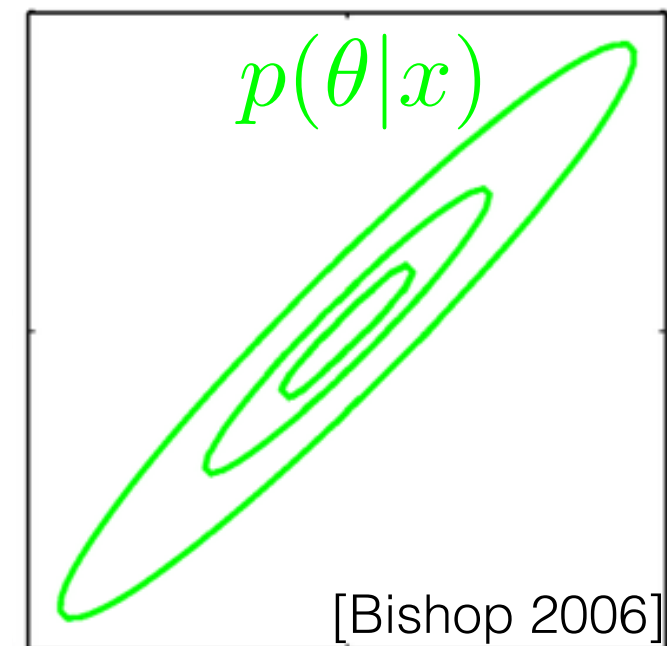
- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0}$$





# Linear response

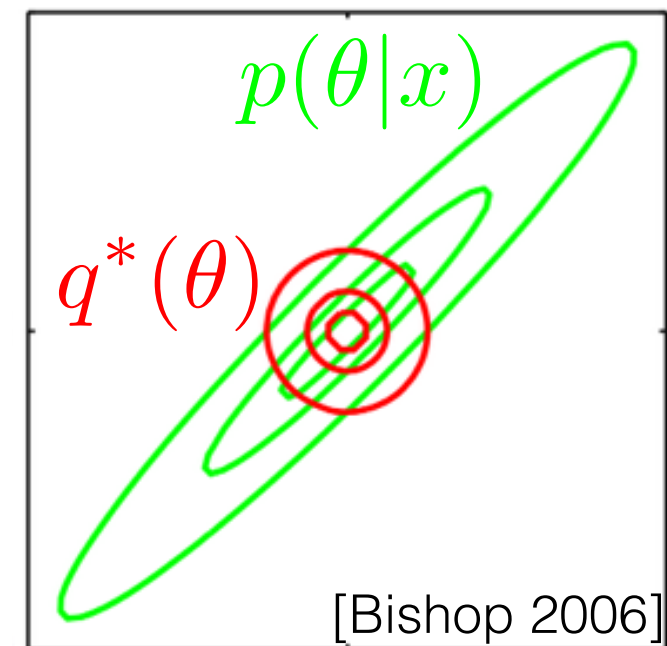
- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0}$$



# Linear response

- Cumulant-generating function

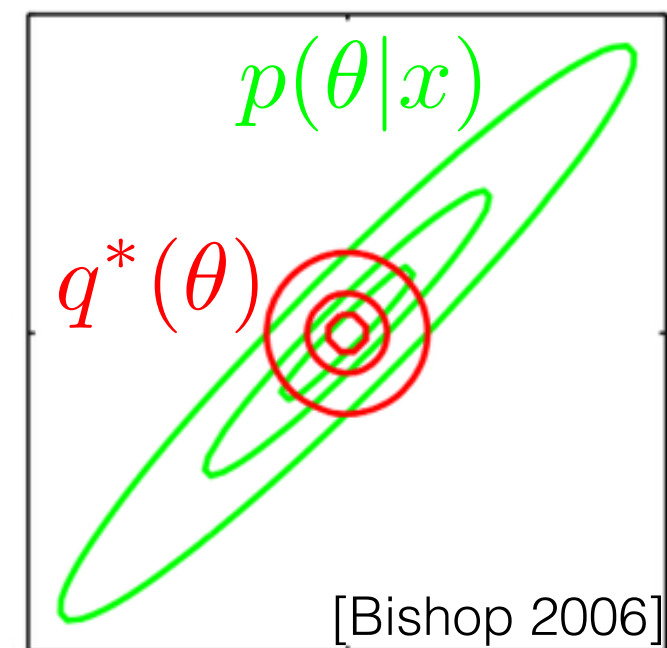
$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0}$$

$$V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$



# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

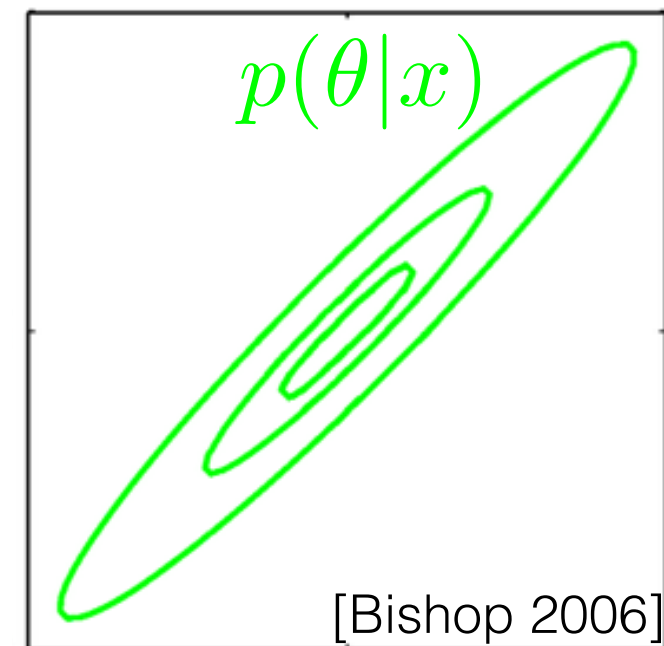
$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0}$$

$$V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”



# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

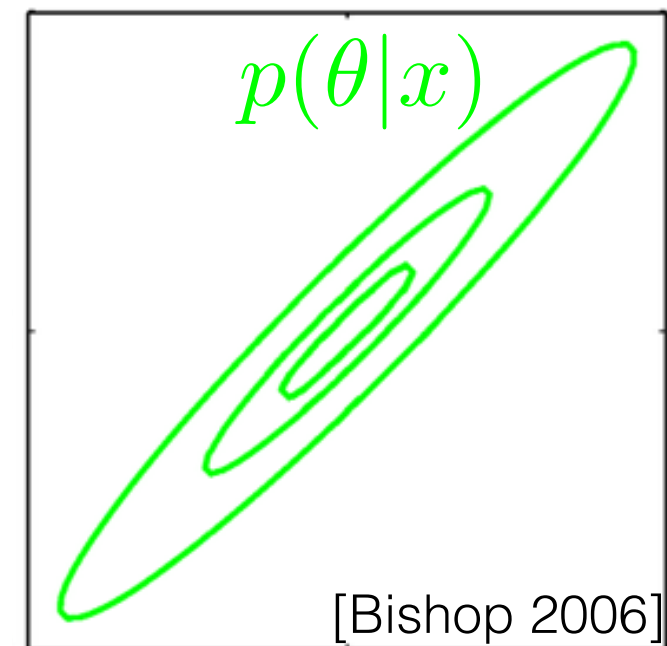
- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0}$$

$$V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p(\theta|x)$$



# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

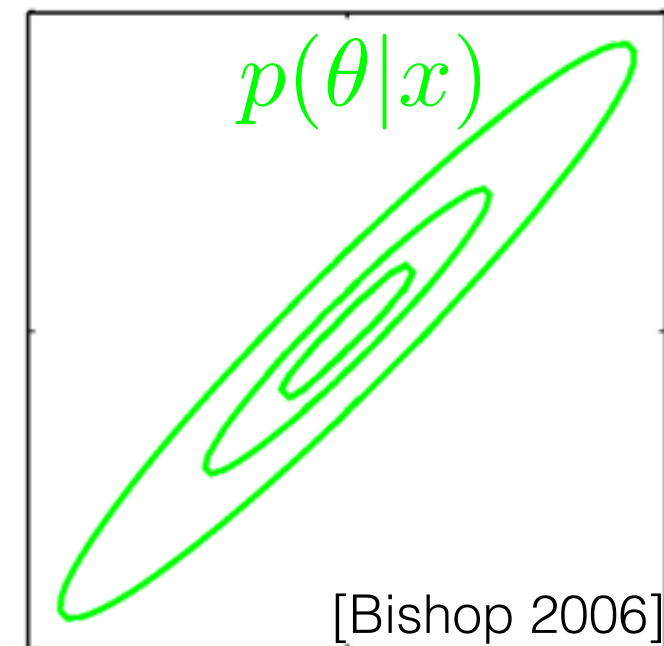
- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0}$$

$$V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p(\theta|x) + t^T \theta$$



# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

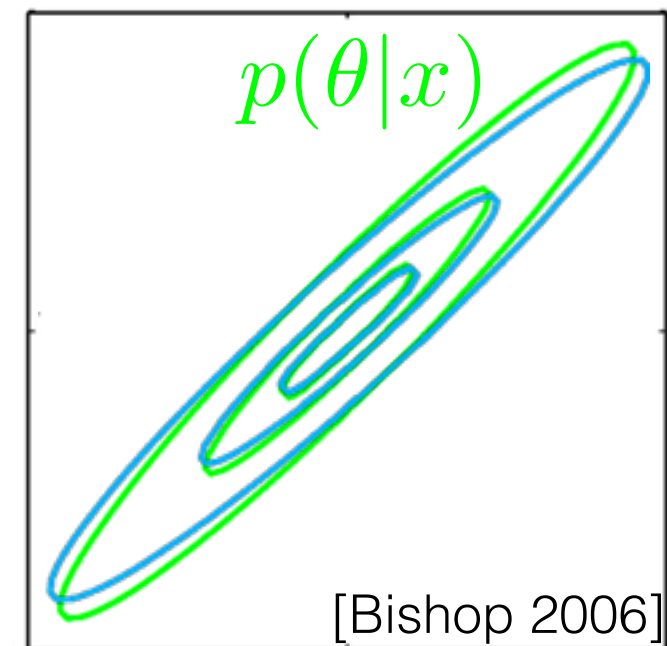
- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0}$$

$$V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p(\theta|x) + t^T \theta$$



# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

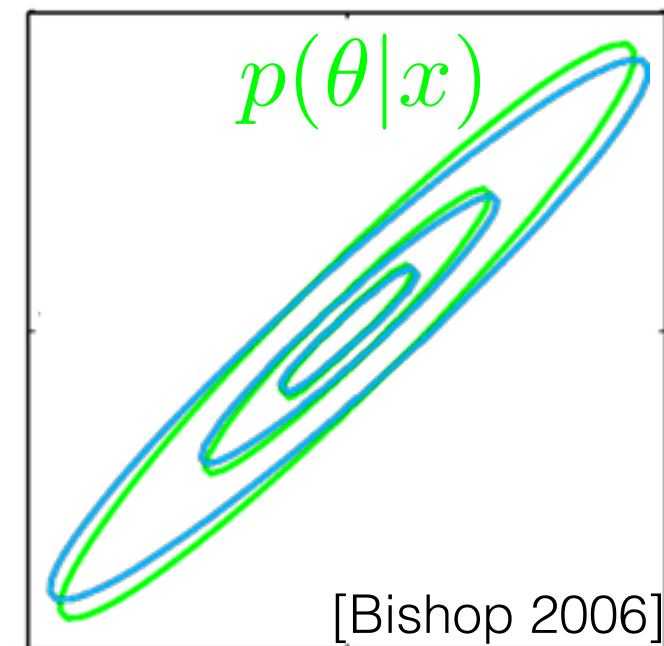
- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0}$$

$$V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta$$



# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta}$$

$$\text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

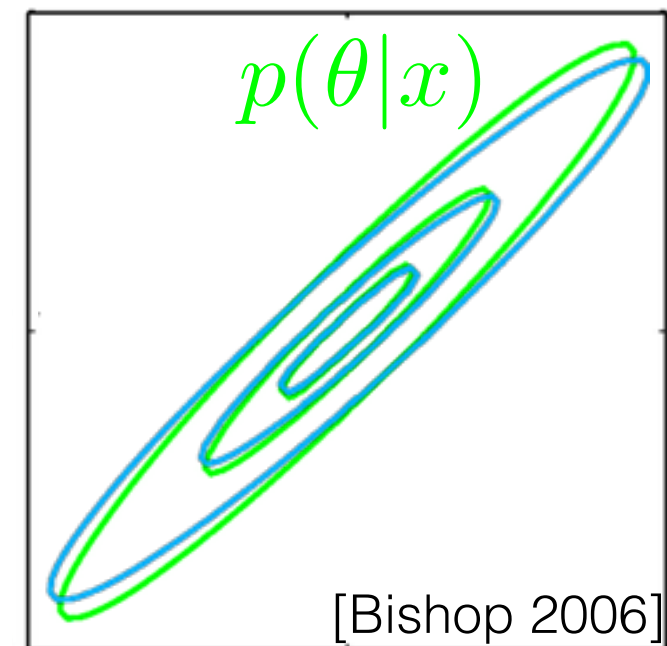
- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0}$$

$$V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t)$$





# Linear response

- Cumulant-generating function

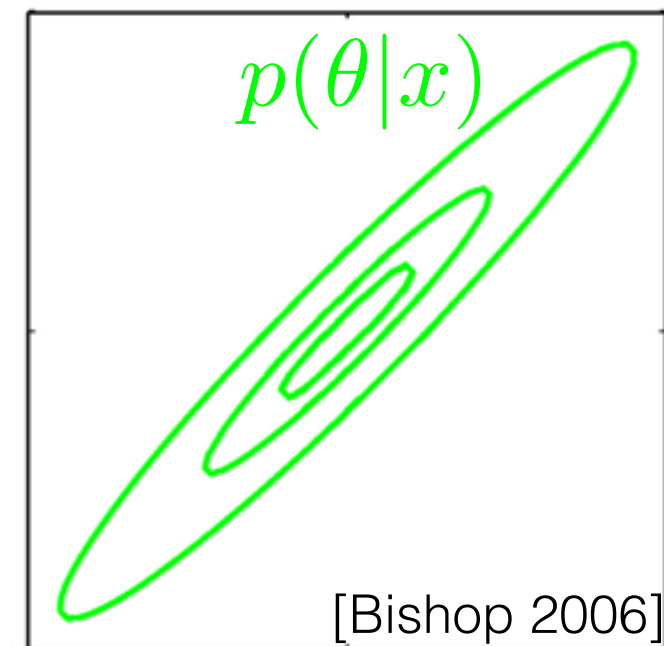
$$C(t) := \log \mathbb{E} e^{t^T \theta} \quad \text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0} \quad V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t)$$



# Linear response

- Cumulant-generating function

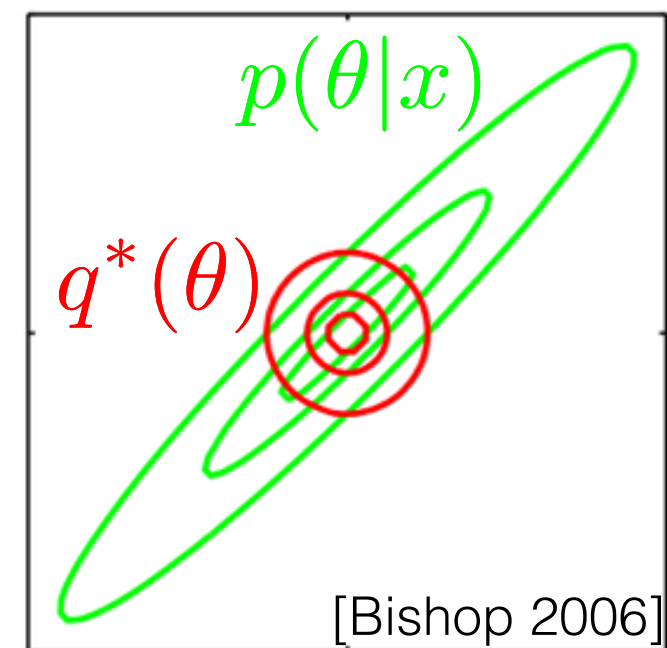
$$C(t) := \log \mathbb{E} e^{t^T \theta} \quad \text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0} \quad V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t), \text{ MFVB } q_t^*$$



# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta} \quad \text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

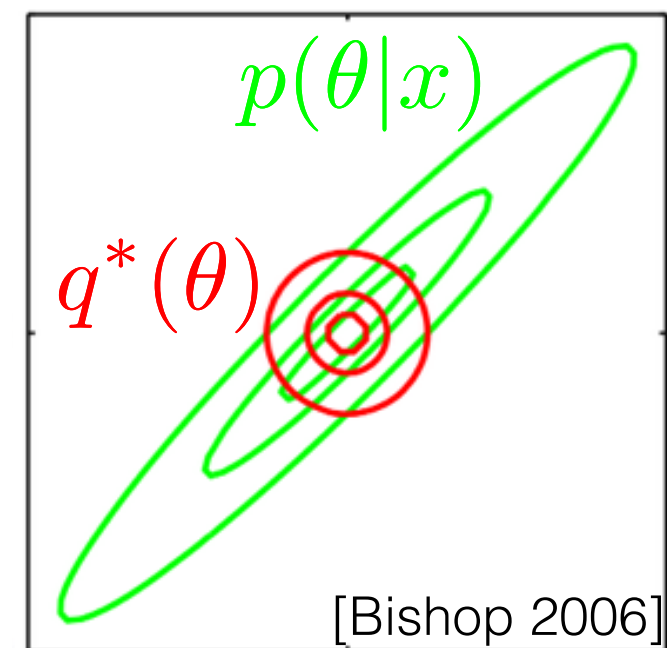
- Exact posterior covariance vs MFVB covariance

$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0} \quad V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t), \text{ MFVB } q_t^*$$

- The LRVB approximation



# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta} \quad \text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

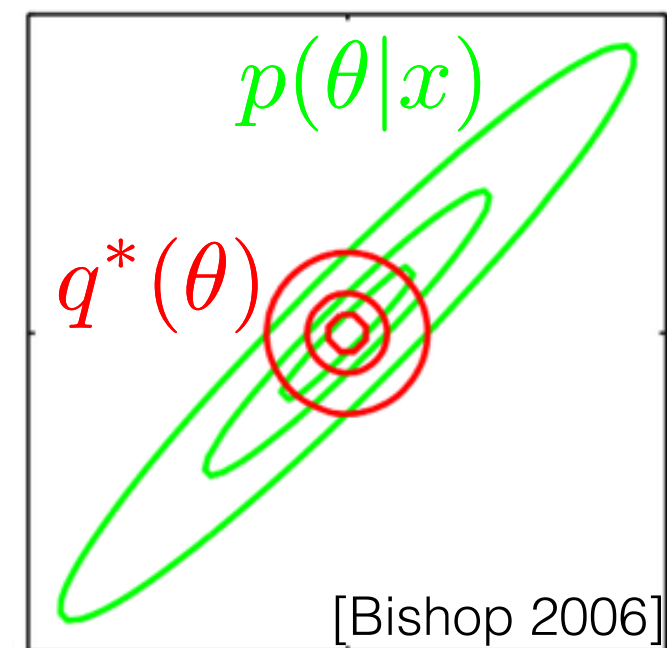
$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0} \quad V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t), \text{ MFVB } q_t^*$$

- The LRVB approximation

$$\Sigma = \frac{d}{dt^T} \left[ \frac{d}{dt} C_{p(\cdot|x)}(t) \right] \Big|_{t=0}$$



[see also Oppen, Winther 2003]

# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta} \quad \text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

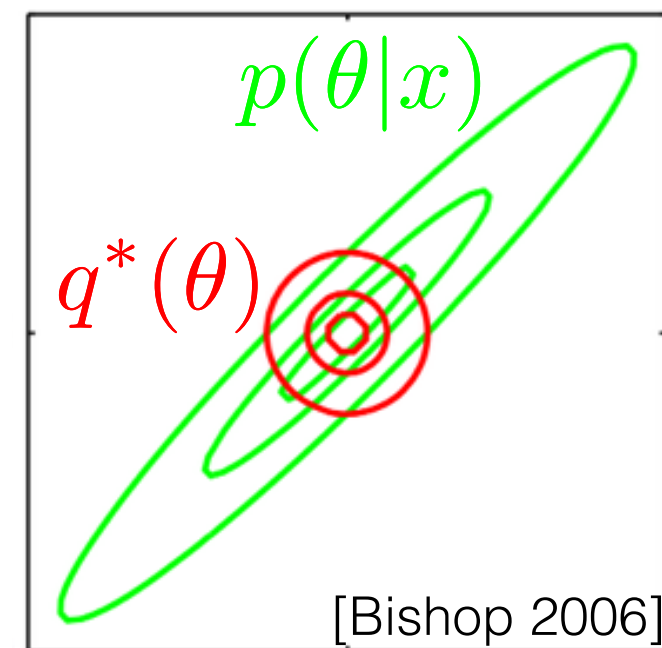
$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0} \quad V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t), \text{ MFVB } q_t^*$$

- The LRVB approximation

$$\Sigma = \left. \frac{d}{dt^T} \mathbb{E}_{p_t} \theta \right|_{t=0}$$



[see also Oppen, Winther 2003]

# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta} \quad \text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

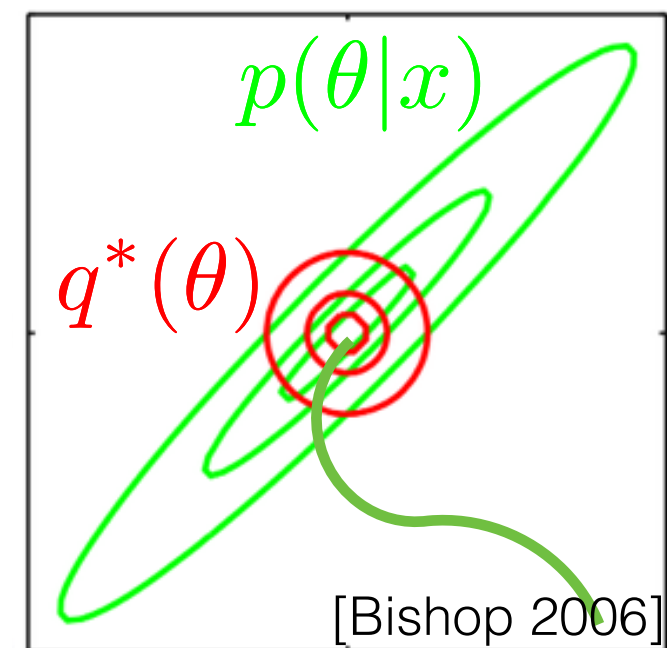
$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0} \quad V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t), \text{ MFVB } q_t^*$$

- The LRVB approximation

$$\Sigma = \left. \frac{d}{dt^T} \mathbb{E}_{p_t} \theta \right|_{t=0}$$



[see also Opper, Winther 2003]

# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta} \quad \text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

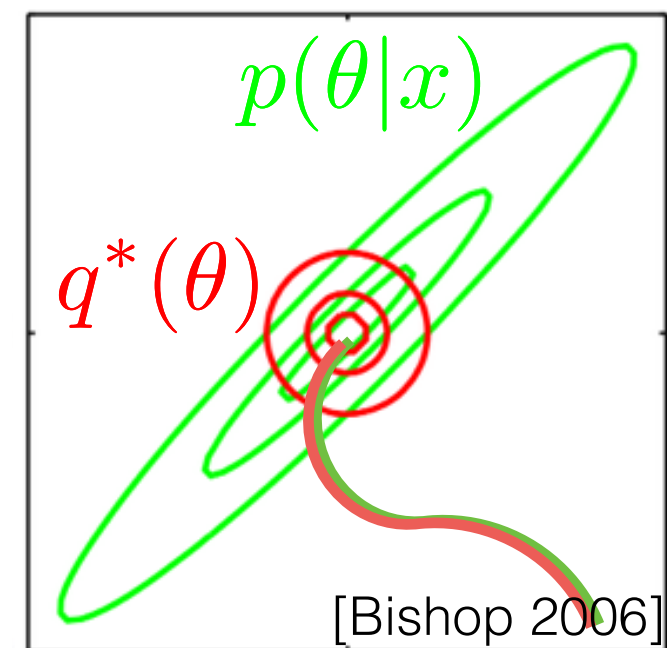
$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0} \quad V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t), \text{ MFVB } q_t^*$$

- The LRVB approximation

$$\Sigma = \left. \frac{d}{dt^T} \mathbb{E}_{p_t} \theta \right|_{t=0}$$



[see also Oppen, Winther 2003]

# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta} \quad \text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

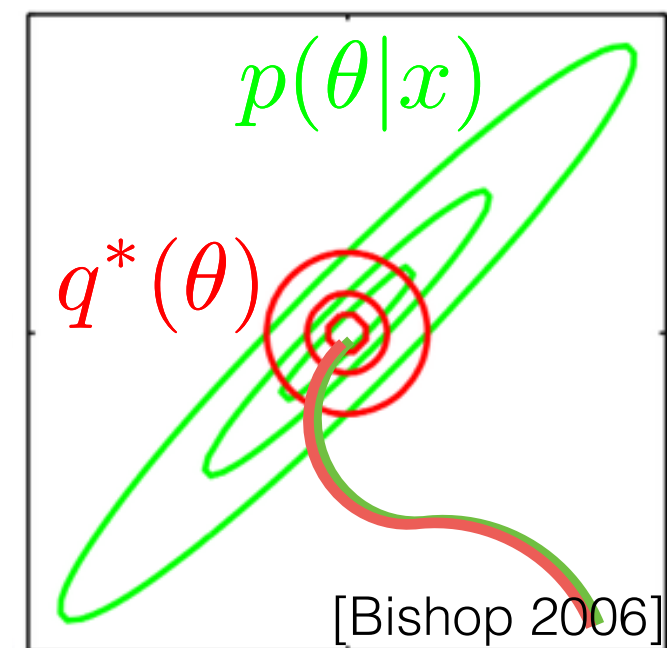
$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0} \quad V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t), \text{ MFVB } q_t^*$$

- The LRVB approximation

$$\Sigma = \left. \frac{d}{dt^T} \mathbb{E}_{p_t} \theta \right|_{t=0} \approx \left. \frac{d}{dt^T} \mathbb{E}_{q_t^*} \theta \right|_{t=0}$$



[see also Opper, Winther 2003]



# Linear response

- Cumulant-generating function

$$C(t) := \log \mathbb{E} e^{t^T \theta} \quad \text{mean} = \left. \frac{d}{dt} C(t) \right|_{t=0}$$

- Exact posterior covariance vs MFVB covariance

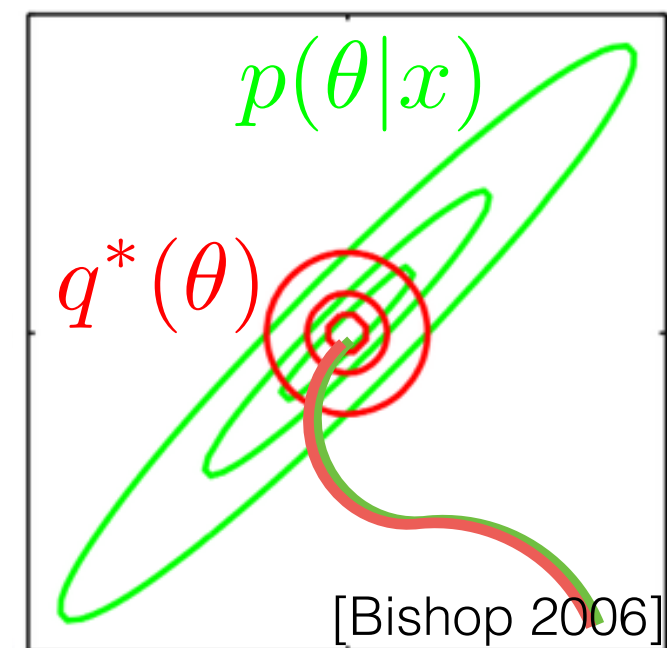
$$\Sigma := \left. \frac{d^2}{dt^T dt} C_{p(\cdot|x)}(t) \right|_{t=0} \quad V := \left. \frac{d^2}{dt^T dt} C_{q^*}(t) \right|_{t=0}$$

- “Linear response”

$$\log p_t(\theta) := \log p(\theta|x) + t^T \theta - C(t), \text{ MFVB } q_t^*$$

- The LRVB approximation

$$\Sigma = \left. \frac{d}{dt^T} \mathbb{E}_{p_t} \theta \right|_{t=0} \approx \left. \frac{d}{dt^T} \mathbb{E}_{q_t^*} \theta \right|_{t=0} =: \hat{\Sigma}$$



[see also Opper, Winther 2003]

# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \left. \frac{d}{dt} \mathbb{E}_{q_t^*} \theta \right|_{t=0}$

# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \frac{d}{dt} \mathbb{E}_{q_t^*} \theta \Big|_{t=0}$
- Suppose  $q_t$  exponential family

# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \frac{d}{dt} \mathbb{E}_{q_t^*} \theta \Big|_{t=0}$
- Suppose  $q_t$  exponential family with mean parametrization  $m_t$

# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \left. \frac{d}{dt} \mathbb{E}_{q_t^*} \theta \right|_{t=0}$
- Suppose  $q_t$  exponential family with mean parametrization  $m_t$

$$\hat{\Sigma} =$$

# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \frac{d}{dt^T} \mathbb{E}_{q_t^*} \theta \Big|_{t=0}$
- Suppose  $q_t$  exponential family with mean parametrization  $m_t$

$$\hat{\Sigma} = \left( \frac{\partial^2 KL}{\partial m \partial m^T} \Big|_{m=m^*} \right)^{-1}$$

# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \frac{d}{dt^T} \mathbb{E}_{q_t^*} \theta \Big|_{t=0}$
- Suppose  $q_t$  exponential family with mean parametrization  $m_t$

$$\hat{\Sigma} = \left( \frac{\partial^2 KL}{\partial m \partial m^T} \Big|_{m=m^*} \right)^{-1} = (I - VH)^{-1}V$$

# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \frac{d}{dt^T} \mathbb{E}_{q_t^*} \theta \Big|_{t=0}$
- Suppose  $q_t$  exponential family with mean parametrization  $m_t$

$$\hat{\Sigma} = \left( \frac{\partial^2 KL}{\partial m \partial m^T} \Big|_{m=m^*} \right)^{-1} = (I - VH)^{-1}V$$

- Symmetric and positive definite at local min of KL



# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \frac{d}{dt} \mathbb{E}_{q_t^*} \theta \Big|_{t=0}$

- Suppose  $q_t$  exponential family with mean parametrization  $m_t$

$$\hat{\Sigma} = \left( \frac{\partial^2 KL}{\partial m \partial m^T} \Big|_{m=m^*} \right)^{-1} = (I - VH)^{-1}V$$

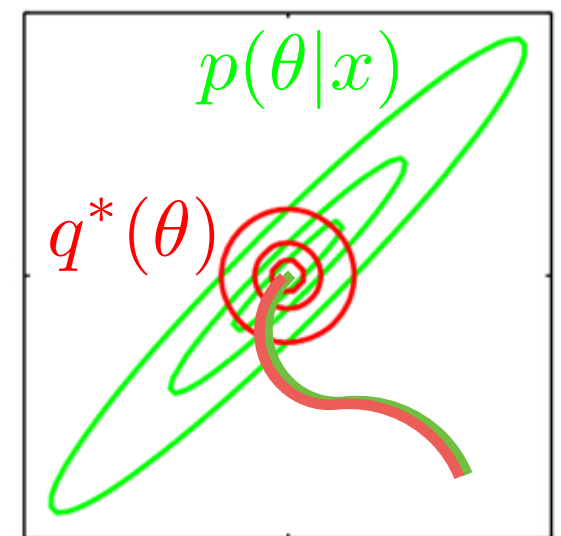
- Symmetric and positive definite at local min of KL
- The LRVB assumption:  $\mathbb{E}_{p_t} \theta \approx \mathbb{E}_{q_t^*} \theta$

# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \frac{d}{dt} \mathbb{E}_{q_t^*} \theta \Big|_{t=0}$
- Suppose  $q_t$  exponential family with mean parametrization  $m_t$

$$\hat{\Sigma} = \left( \frac{\partial^2 KL}{\partial m \partial m^T} \Big|_{m=m^*} \right)^{-1} = (I - VH)^{-1}V$$

- Symmetric and positive definite at local min of KL
- The LRVB assumption:  $\mathbb{E}_{p_t} \theta \approx \mathbb{E}_{q_t^*} \theta$



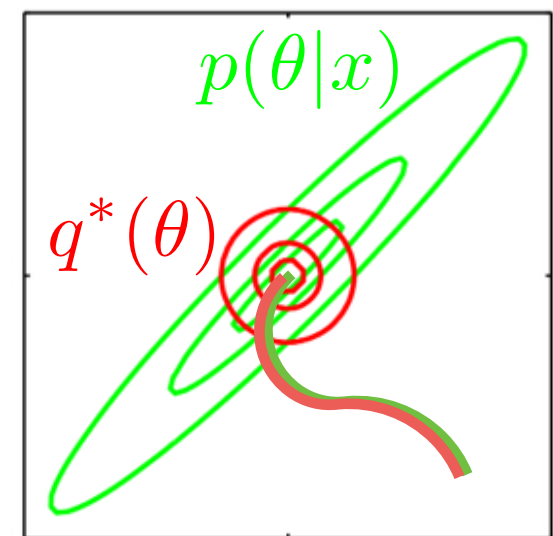
[Bishop 2006]

# LRVB estimator

- LRVB covariance estimate  $\hat{\Sigma} := \left. \frac{d}{dt} \mathbb{E}_{q_t^*} \theta \right|_{t=0}$
- Suppose  $q_t$  exponential family with mean parametrization  $m_t$

$$\hat{\Sigma} = \left( \left. \frac{\partial^2 KL}{\partial m \partial m^T} \right|_{m=m^*} \right)^{-1} = (I - VH)^{-1}V$$

- Symmetric and positive definite at local min of KL
- The LRVB assumption:  $\mathbb{E}_{p_t} \theta \approx \mathbb{E}_{q_t^*} \theta$
- LRVB estimate is exact when MFVB gives exact mean (e.g. multivariate normal)



[Bishop 2006]

# Microcredit Experiment

# Microcredit Experiment

- Simplified from Meager (2016)

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )


# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:



# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit   $y_{kn}$


# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit  $\rightarrow y_{kn} \stackrel{\text{indep}}{\sim} \mathcal{N}(\quad, \quad)$

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit   $y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k, \sigma_k^2)$

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit 

$$y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn}\tau_k, \quad )$$

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit  $\rightarrow y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn}\tau_k, \quad )$

$T_{kn}$  is highlighted in green.

1 if microcredit

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit  $\rightarrow y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn} \tau_k, \quad )$

$\swarrow$  1 if microcredit

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit  $\rightarrow y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn}\tau_k, \sigma_k^2)$

$\swarrow$  1 if microcredit

# Microcredit Experiment

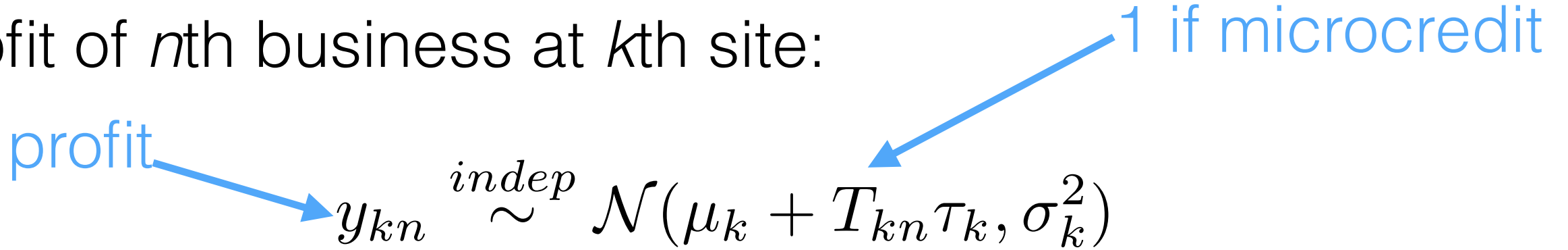
- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit  $\rightarrow y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn}\tau_k, \sigma_k^2)$

$\swarrow$  1 if microcredit



# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:  

$$y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn}\tau_k, \sigma_k^2)$$
- Priors and hyperpriors:

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit  $\rightarrow y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn}\tau_k, \sigma_k^2)$

$\swarrow$  1 if microcredit

- Priors and hyperpriors:

$$\begin{pmatrix} \mu_k \\ \tau_k \end{pmatrix} \stackrel{iid}{\sim} \mathcal{N}\left(\begin{pmatrix} \mu \\ \tau \end{pmatrix}, C\right)$$

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit  $\rightarrow y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn}\tau_k, \sigma_k^2)$

1 if microcredit  $\rightarrow T_{kn}$

- Priors and hyperpriors:

$$\begin{pmatrix} \mu_k \\ \tau_k \end{pmatrix} \stackrel{iid}{\sim} \mathcal{N}\left(\begin{pmatrix} \mu \\ \tau \end{pmatrix}, C\right)$$

$$\sigma_k^{-2} \stackrel{iid}{\sim} \Gamma(a, b)$$

# Microcredit Experiment

- Simplified from Meager (2016)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit  $\rightarrow y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn}\tau_k, \sigma_k^2)$

1 if microcredit  $\rightarrow T_{kn}$

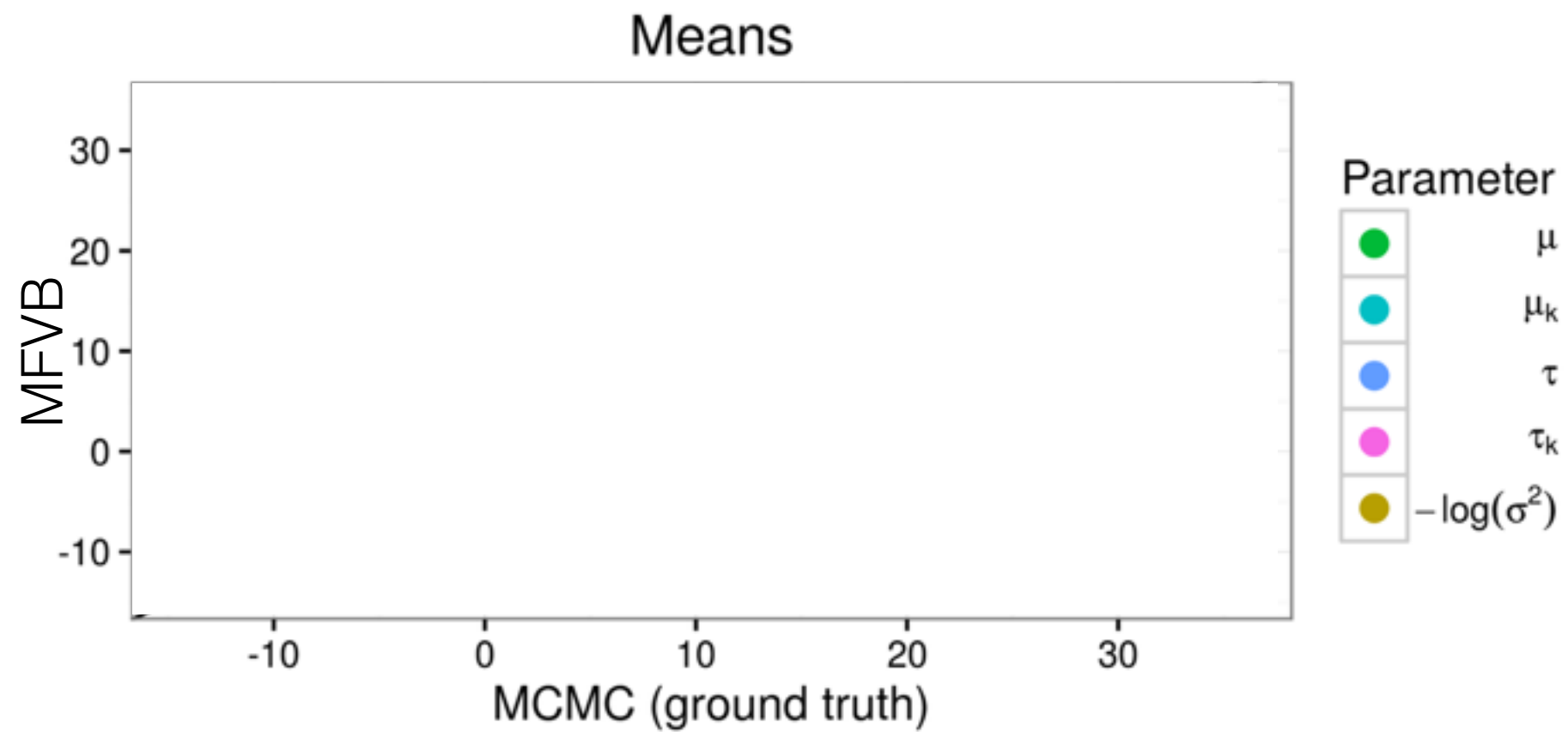
- Priors and hyperpriors:

$$\begin{pmatrix} \mu_k \\ \tau_k \end{pmatrix} \stackrel{iid}{\sim} \mathcal{N}\left(\begin{pmatrix} \mu \\ \tau \end{pmatrix}, C\right) \quad \begin{pmatrix} \mu \\ \tau \end{pmatrix} \stackrel{iid}{\sim} \mathcal{N}\left(\begin{pmatrix} \mu_0 \\ \tau_0 \end{pmatrix}, \Lambda^{-1}\right)$$

$$\sigma_k^{-2} \stackrel{iid}{\sim} \Gamma(a, b)$$

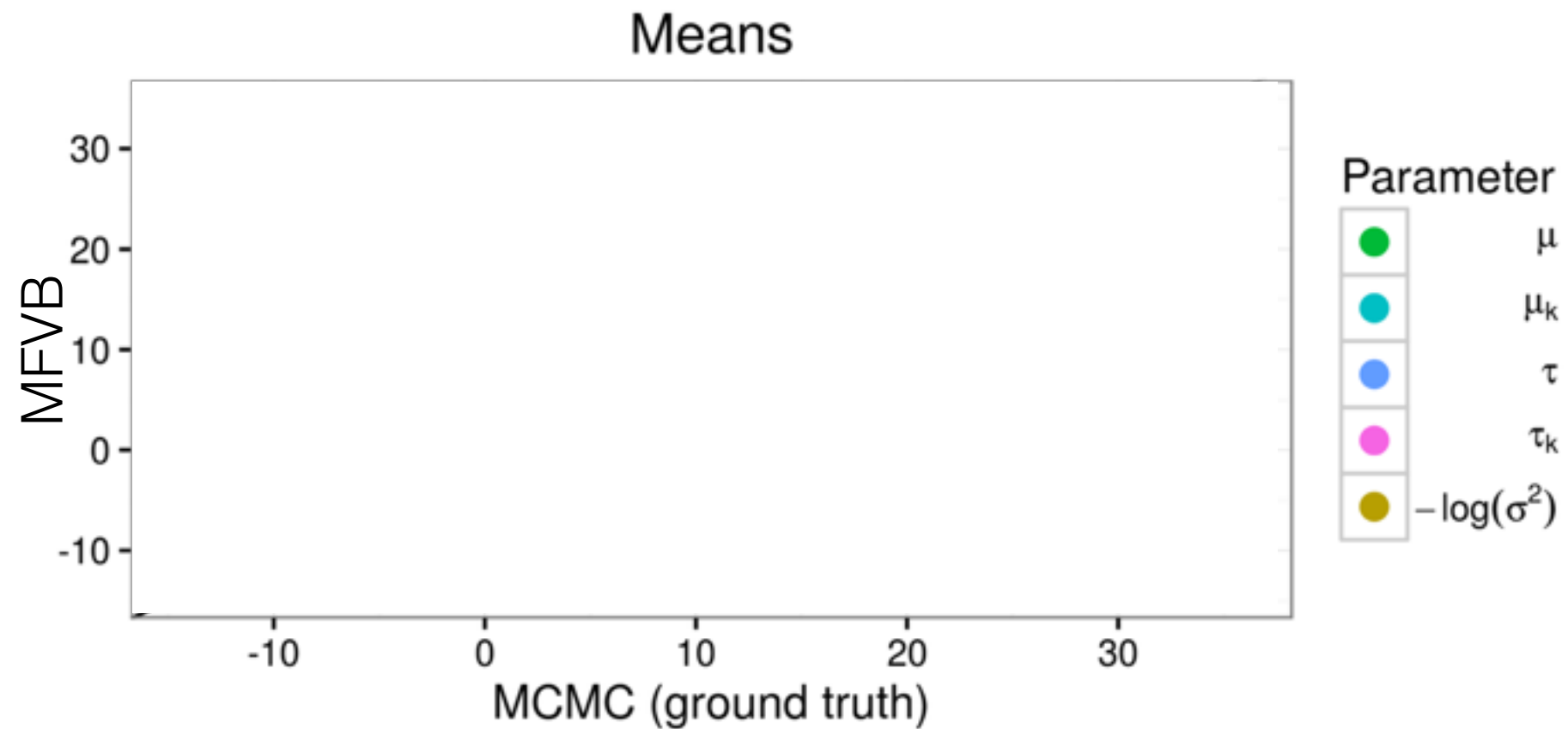
$$C \sim \text{Sep\&LKJ}(\eta, c, d)$$

# Microcredit Experiment



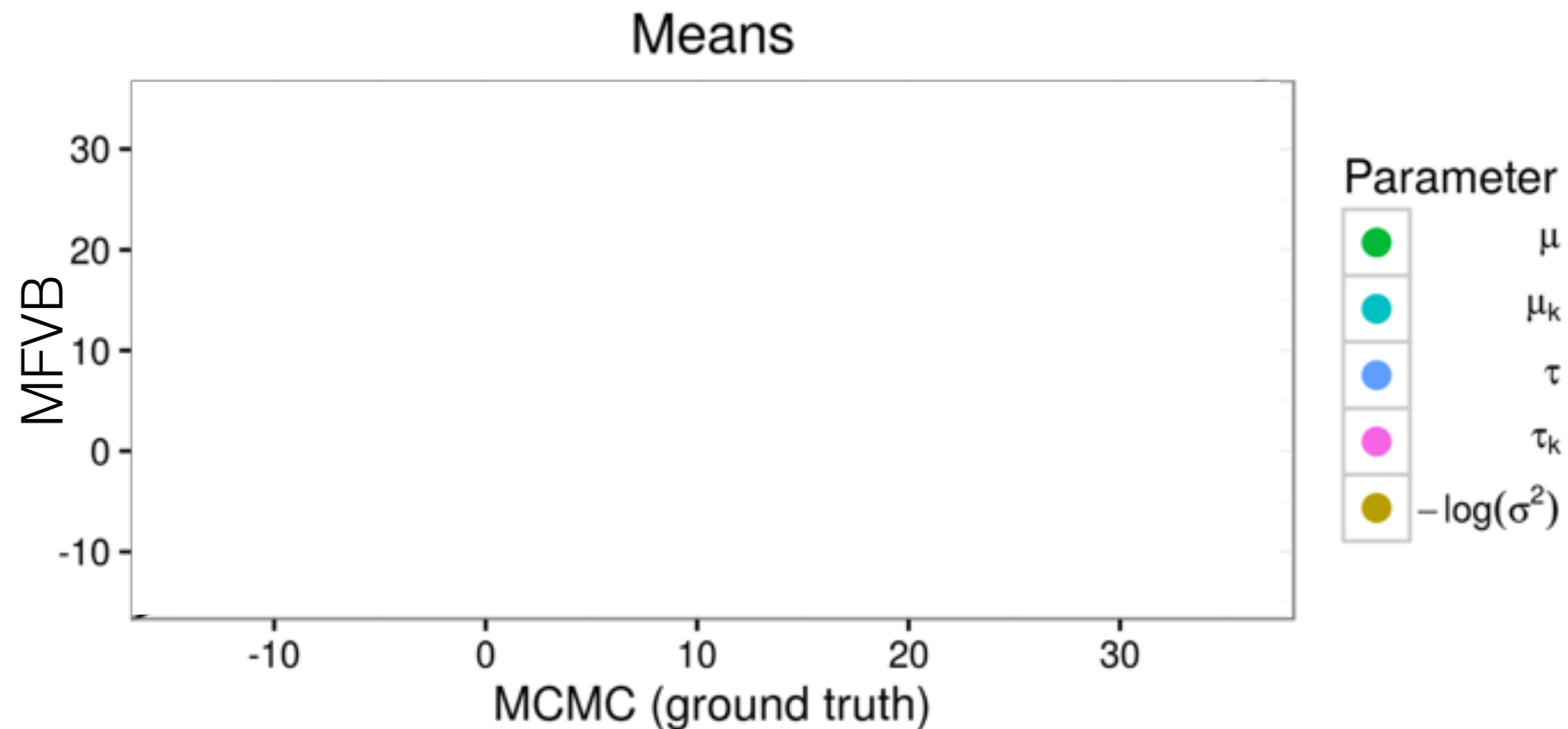
# Microcredit Experiment

- *One set* of 2500 MCMC draws:  
**45 minutes**



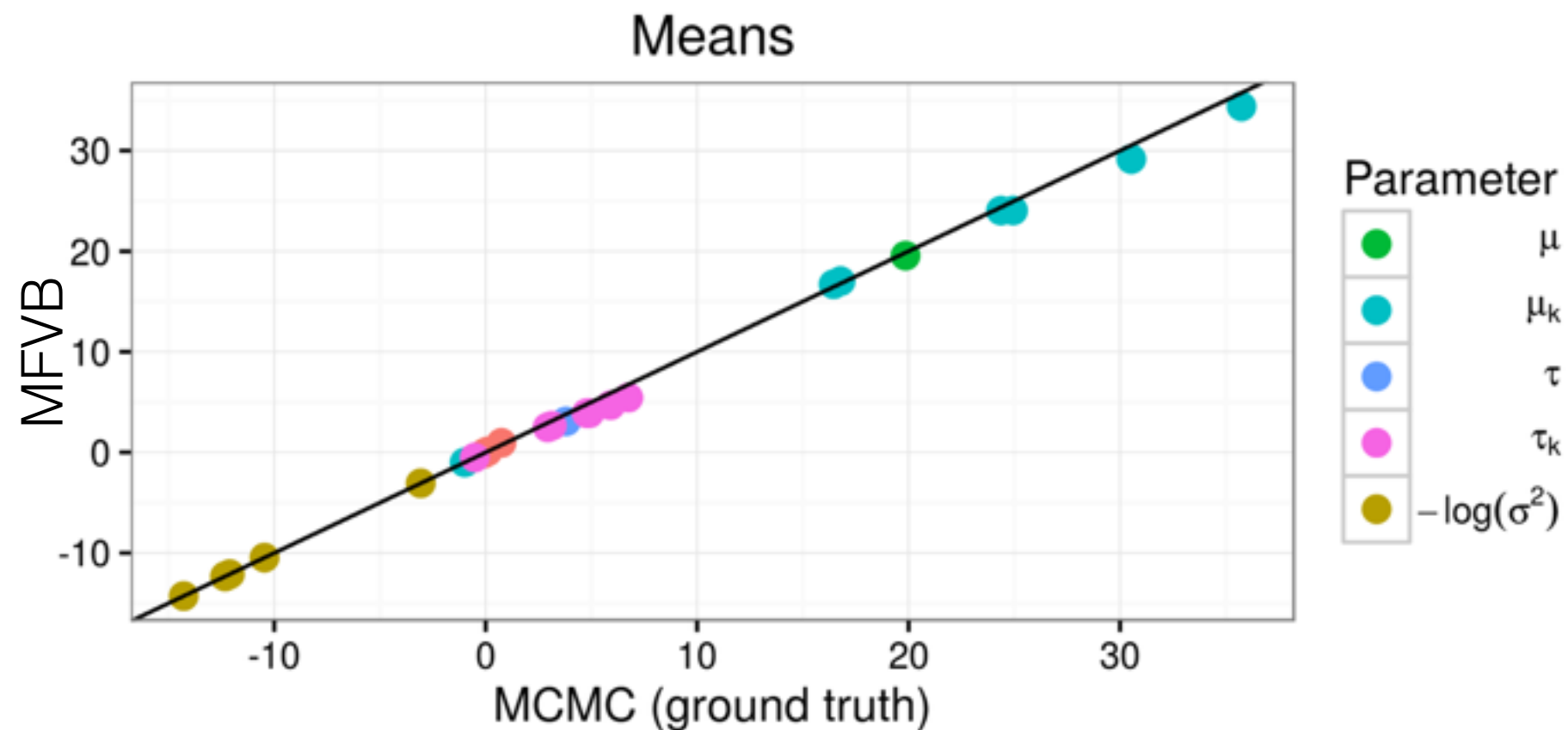
# Microcredit Experiment

- *One set of 2500* MCMC draws:  
**45 minutes**
- All of MFVB optimization, LRVB uncertainties, all sensitivity measures:  
**58 seconds**



# Microcredit Experiment

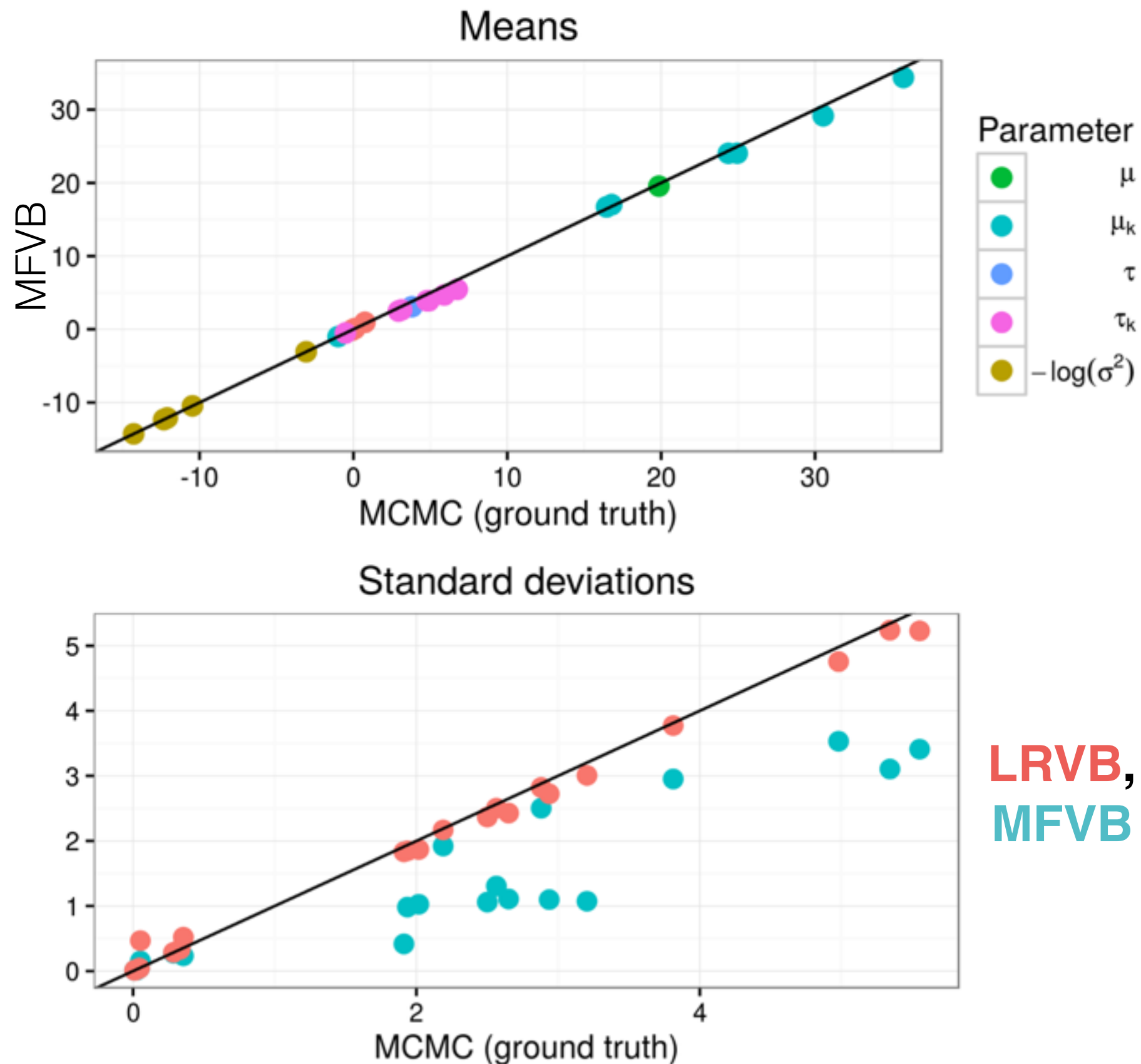
- *One set of 2500* MCMC draws:  
**45 minutes**
- All of MFVB optimization, LRVB uncertainties, all sensitivity measures:  
**58 seconds**





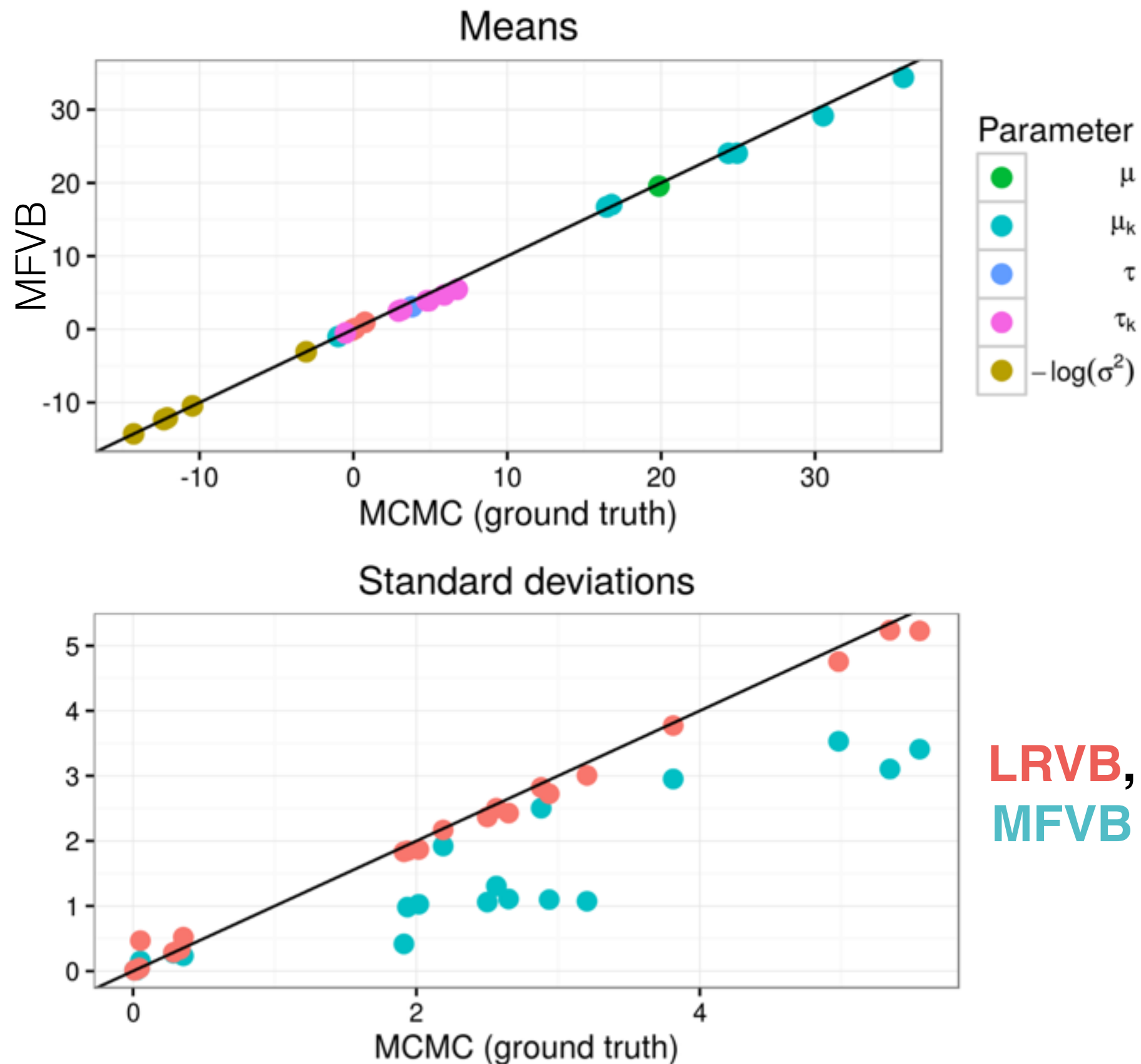
# Microcredit Experiment

- One set of 2500 MCMC draws:  
**45 minutes**
- All of MFVB optimization, LRVB uncertainties, all sensitivity measures:  
**58 seconds**



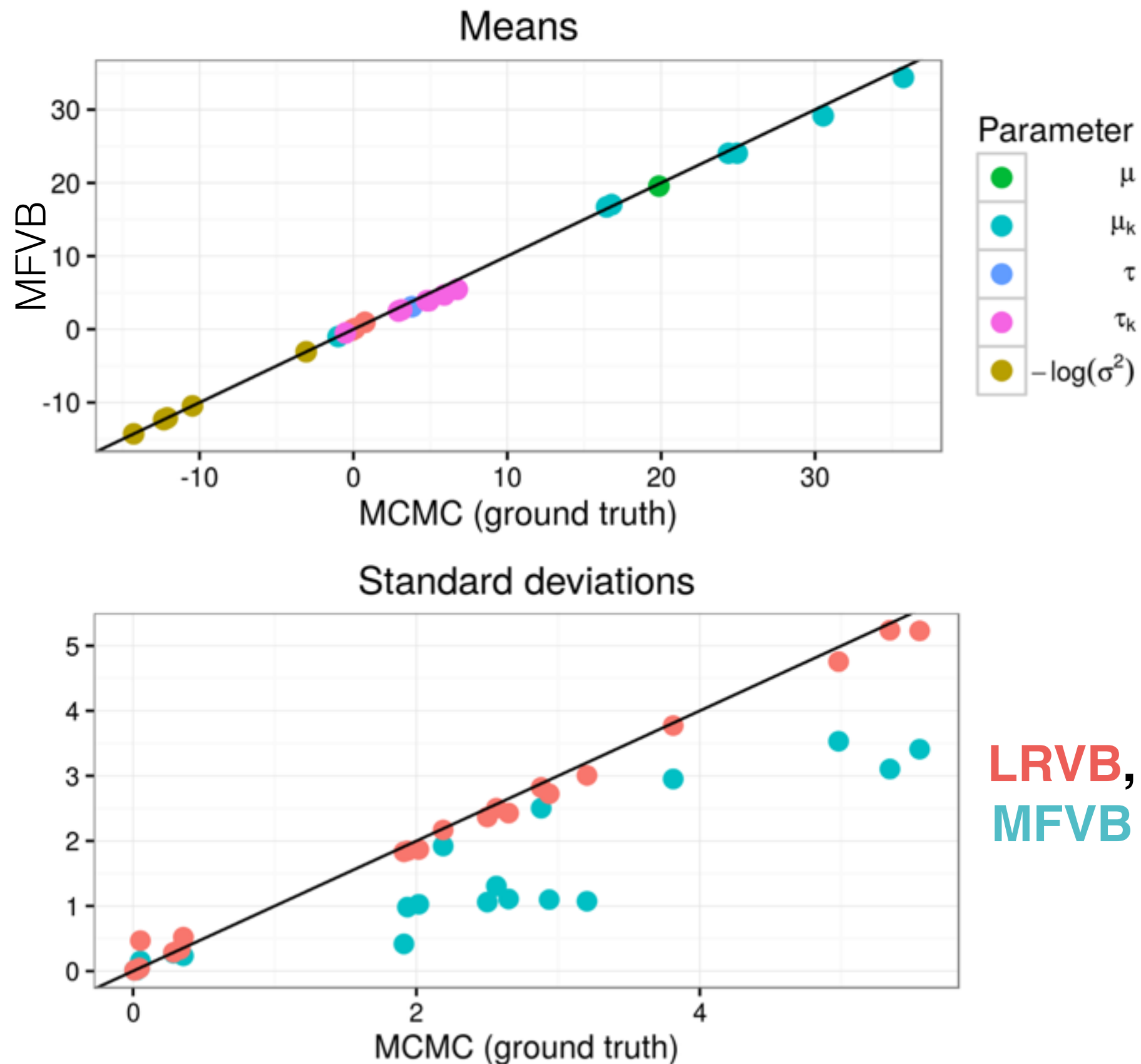
# Microcredit Experiment

- One set of 2500 MCMC draws:  
**45 minutes**
- All of MFVB optimization, LRVB uncertainties, all sensitivity measures:  
**58 seconds**
- $\tau$  mean (MFVB):  
3.08 USD PPP



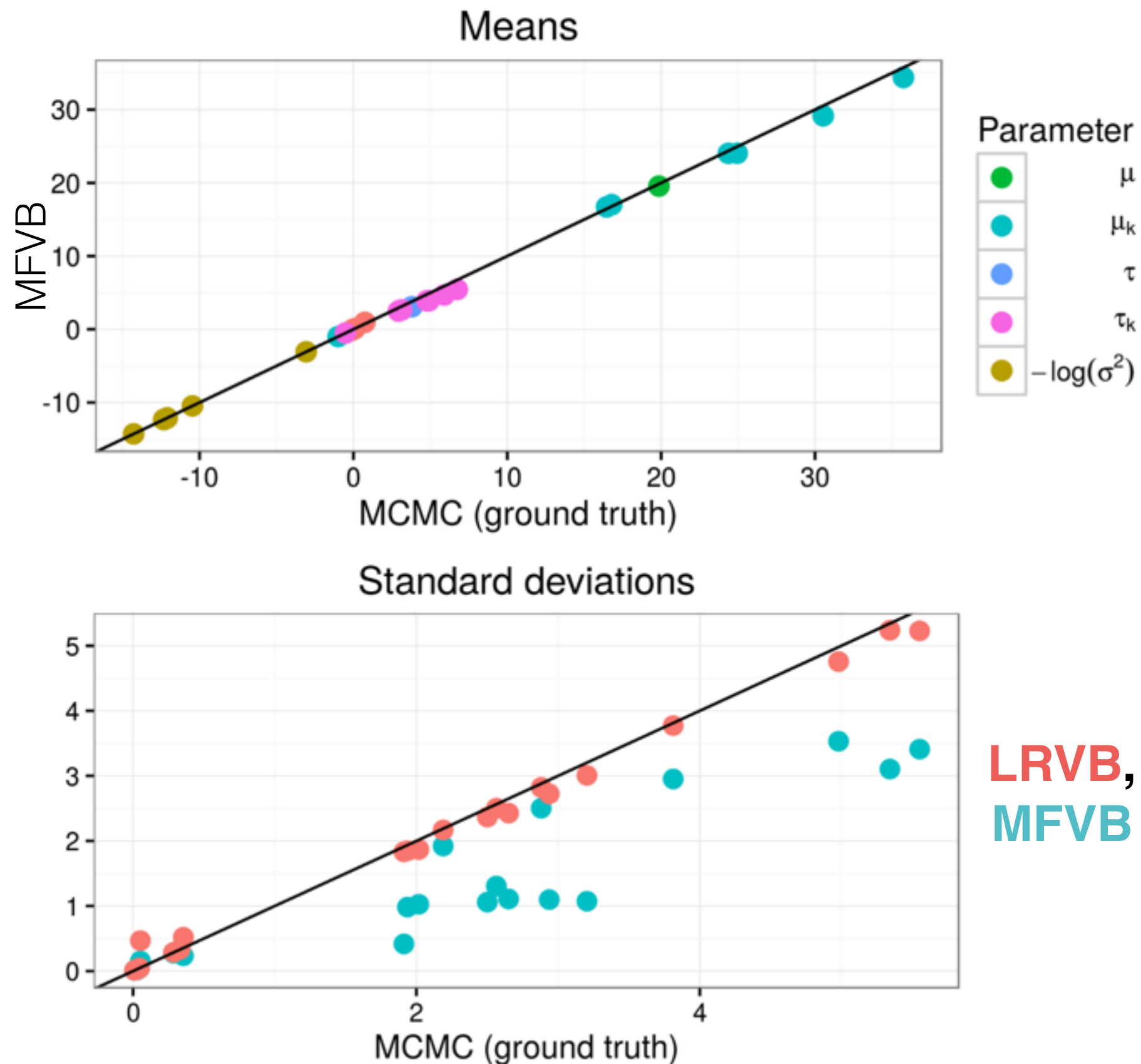
# Microcredit Experiment

- One set of 2500 MCMC draws:  
**45 minutes**
- All of MFVB optimization, LRVB uncertainties, all sensitivity measures:  
**58 seconds**
- $\tau$  mean (MFVB):  
3.08 USD PPP
- $\tau$  std dev (LRVB):  
1.83 USD PPP



# Microcredit Experiment

- One set of 2500 MCMC draws:  
**45 minutes**
- All of MFVB optimization, LRVB uncertainties, all sensitivity measures:  
**58 seconds**
- $\tau$  mean (MFVB):  
3.08 USD PPP
- $\tau$  std dev (LRVB):  
1.83 USD PPP
- Mean is 1.68 std dev from 0



# Experiments

# Experiments

- Gaussian mixture model

$$P(z_{nk} = 1) = \pi_k, \quad p(x|\pi, \mu, \Lambda, z) = \prod_{n=1:N} \prod_{k=1:K} \mathcal{N}(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$$

with conjugate priors on  $\pi, \mu, \Lambda$

# Experiments

- Gaussian mixture model

$$P(z_{nk} = 1) = \pi_k, \quad p(x|\pi, \mu, \Lambda, z) = \prod_{n=1:N} \prod_{k=1:K} \mathcal{N}(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$$

with conjugate priors on  $\pi, \mu, \Lambda$

- 68 simulated data sets (2 components, 2 dimensions),  
10,000 data points each, R `bayesm` package

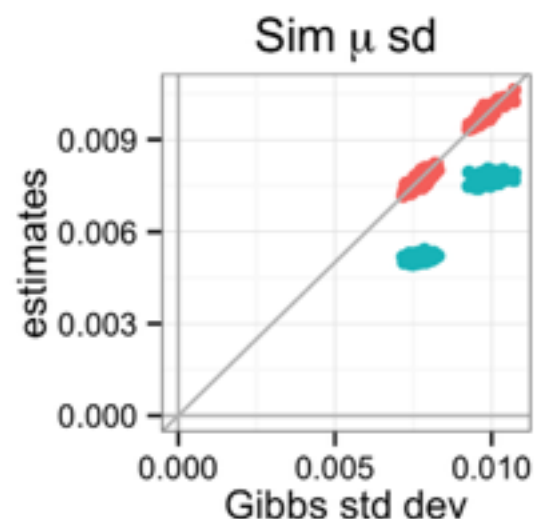
# Experiments

- Gaussian mixture model

$$P(z_{nk} = 1) = \pi_k, \quad p(x|\pi, \mu, \Lambda, z) = \prod_{n=1:N} \prod_{k=1:K} \mathcal{N}(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$$

with conjugate priors on  $\pi, \mu, \Lambda$

- 68 simulated data sets (2 components, 2 dimensions), 10,000 data points each, R `bayesm` package



**LRVB**,  
**MFVB**



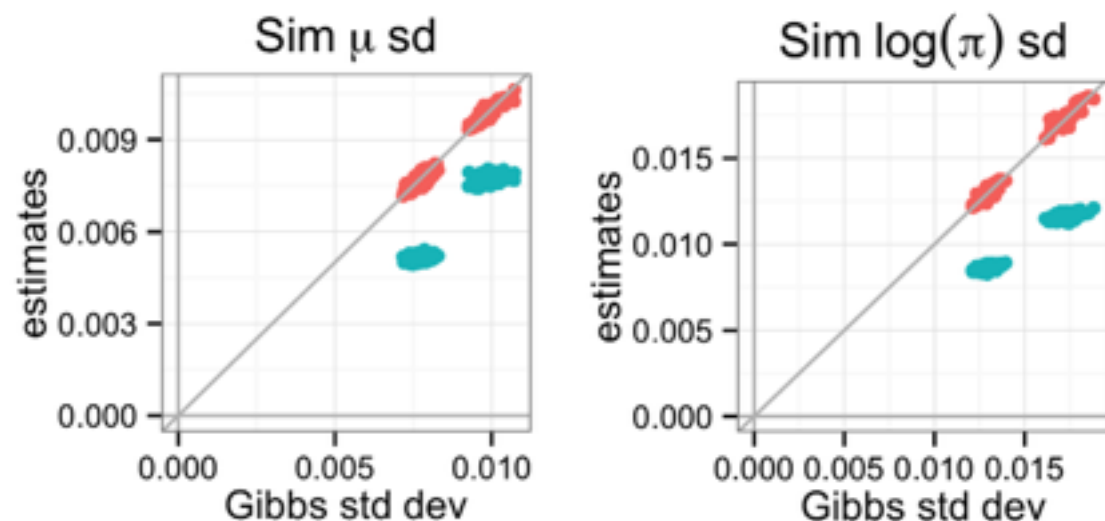
# Experiments

- Gaussian mixture model

$$P(z_{nk} = 1) = \pi_k, \quad p(x|\pi, \mu, \Lambda, z) = \prod_{n=1:N} \prod_{k=1:K} \mathcal{N}(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$$

with conjugate priors on  $\pi, \mu, \Lambda$

- 68 simulated data sets (2 components, 2 dimensions), 10,000 data points each, R `bayesm` package



**LRVB**,  
**MFVB**

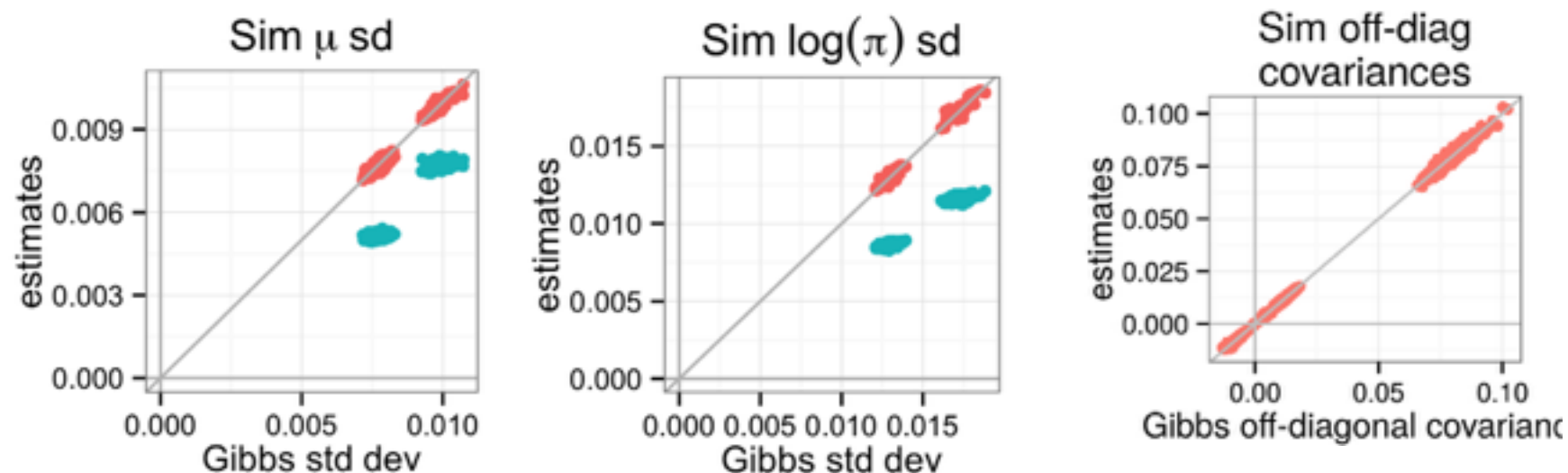
# Experiments

- Gaussian mixture model

$$P(z_{nk} = 1) = \pi_k, \quad p(x|\pi, \mu, \Lambda, z) = \prod_{n=1:N} \prod_{k=1:K} \mathcal{N}(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$$

with conjugate priors on  $\pi, \mu, \Lambda$

- 68 simulated data sets (2 components, 2 dimensions), 10,000 data points each, R `bayesm` package



**LRVB**,  
**MFVB**

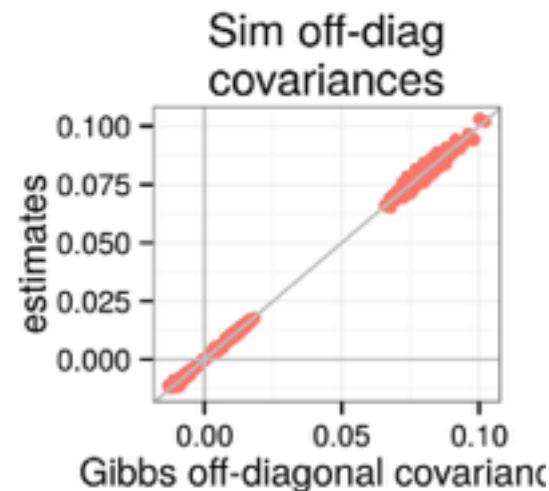
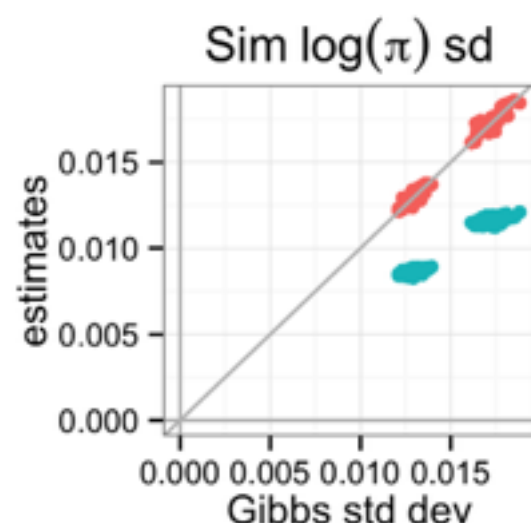
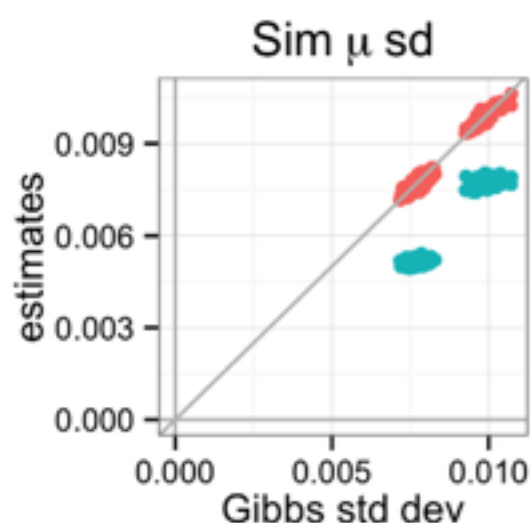
# Experiments

- Gaussian mixture model

$$P(z_{nk} = 1) = \pi_k, \quad p(x|\pi, \mu, \Lambda, z) = \prod_{n=1:N} \prod_{k=1:K} \mathcal{N}(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$$

with conjugate priors on  $\pi, \mu, \Lambda$

- 68 simulated data sets (2 components, 2 dimensions), 10,000 data points each, R `bayesm` package
- MNIST digits: 12,665 0s and 1s; PCA for 25 dimensions



**LRVB**,  
**MFVB**

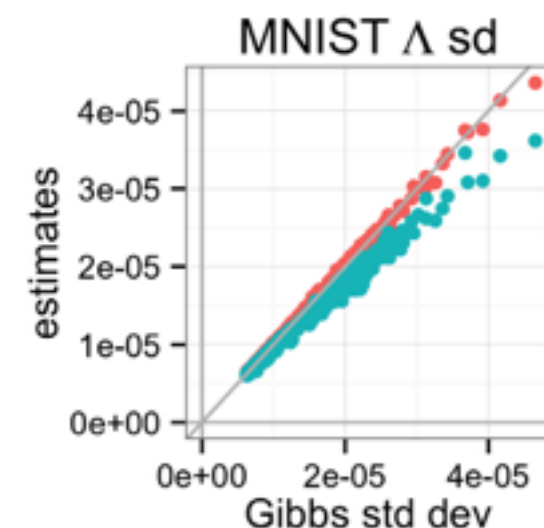
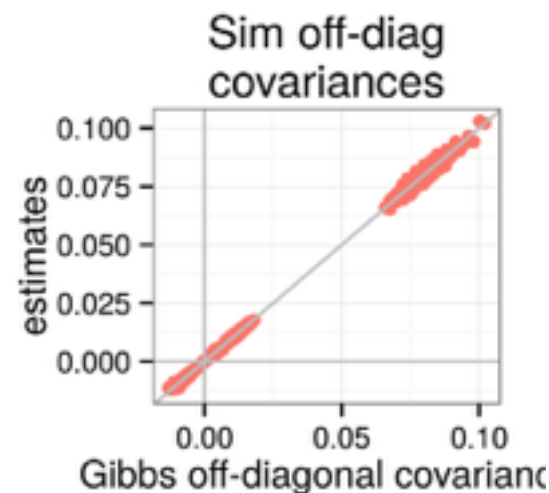
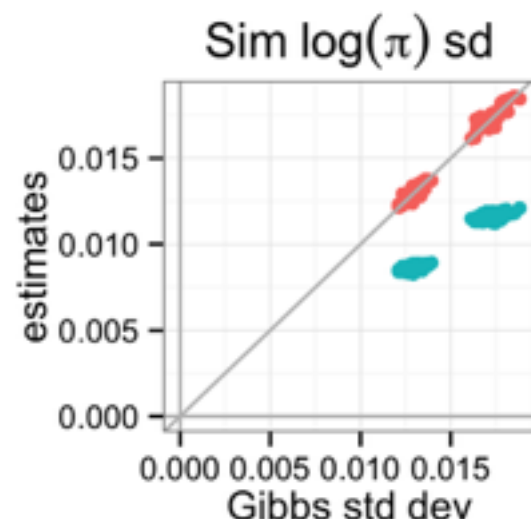
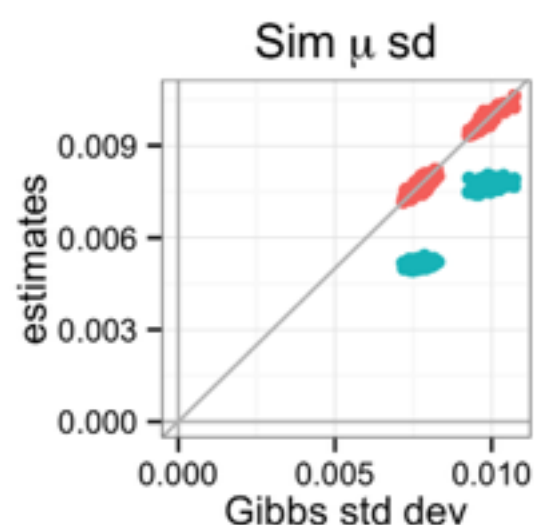
# Experiments

- Gaussian mixture model

$$P(z_{nk} = 1) = \pi_k, \quad p(x|\pi, \mu, \Lambda, z) = \prod_{n=1:N} \prod_{k=1:K} \mathcal{N}(x_n | \mu_k, \Lambda_k^{-1})^{z_{nk}}$$

with conjugate priors on  $\pi, \mu, \Lambda$

- 68 simulated data sets (2 components, 2 dimensions), 10,000 data points each, R `bayesm` package
- MNIST digits: 12,665 0s and 1s; PCA for 25 dimensions



**LRVB**,  
**MFVB**

# Experiments

# Experiments

- Non-conjugate normal-Poisson generalized linear mixed model

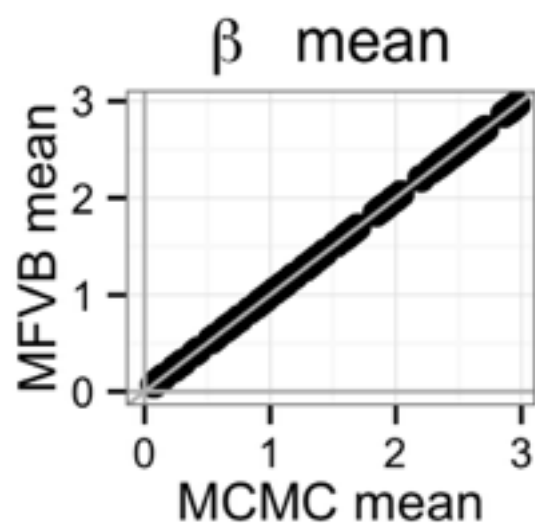
$$z_n | \beta, \tau \stackrel{\text{indep}}{\sim} \mathcal{N}(z_n | \beta x_n, \tau^{-1}), \quad y_n | z_n \stackrel{\text{indep}}{\sim} \text{Poisson}(y_n | \exp(z_n)),$$
$$\beta \sim \mathcal{N}(\beta | 0, \sigma_\beta^2), \quad \tau \sim \text{Gamma}(\tau | \alpha_\tau, \beta_\tau)$$

# Experiments

- Non-conjugate normal-Poisson generalized linear mixed model  
$$z_n | \beta, \tau \stackrel{indep}{\sim} \mathcal{N}(z_n | \beta x_n, \tau^{-1}), \quad y_n | z_n \stackrel{indep}{\sim} \text{Poisson}(y_n | \exp(z_n)),$$
$$\beta \sim \mathcal{N}(\beta | 0, \sigma_\beta^2), \quad \tau \sim \text{Gamma}(\tau | \alpha_\tau, \beta_\tau)$$
- 100 simulated data sets, 500 data points each, R MCMCglmm package

# Experiments

- Non-conjugate normal-Poisson generalized linear mixed model  
$$z_n | \beta, \tau \stackrel{indep}{\sim} \mathcal{N}(z_n | \beta x_n, \tau^{-1}), \quad y_n | z_n \stackrel{indep}{\sim} \text{Poisson}(y_n | \exp(z_n)),$$
$$\beta \sim \mathcal{N}(\beta | 0, \sigma_\beta^2), \quad \tau \sim \text{Gamma}(\tau | \alpha_\tau, \beta_\tau)$$
- 100 simulated data sets, 500 data points each, R MCMCglmm package

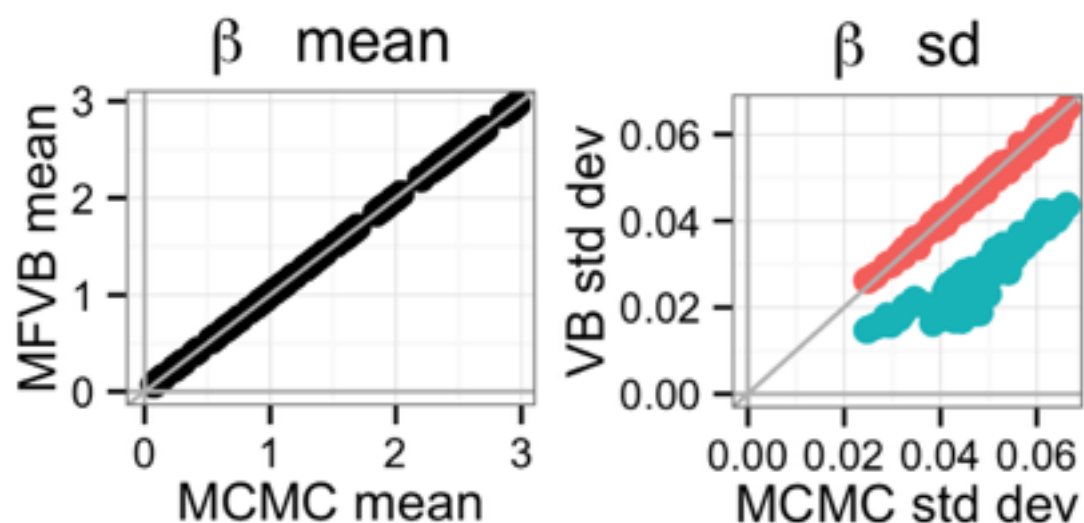




# Experiments

- Non-conjugate normal-Poisson generalized linear mixed model  
$$z_n | \beta, \tau \stackrel{indep}{\sim} \mathcal{N}(z_n | \beta x_n, \tau^{-1}), \quad y_n | z_n \stackrel{indep}{\sim} \text{Poisson}(y_n | \exp(z_n)),$$
$$\beta \sim \mathcal{N}(\beta | 0, \sigma_\beta^2), \quad \tau \sim \text{Gamma}(\tau | \alpha_\tau, \beta_\tau)$$
- 100 simulated data sets, 500 data points each, R  
MCMCglmm package

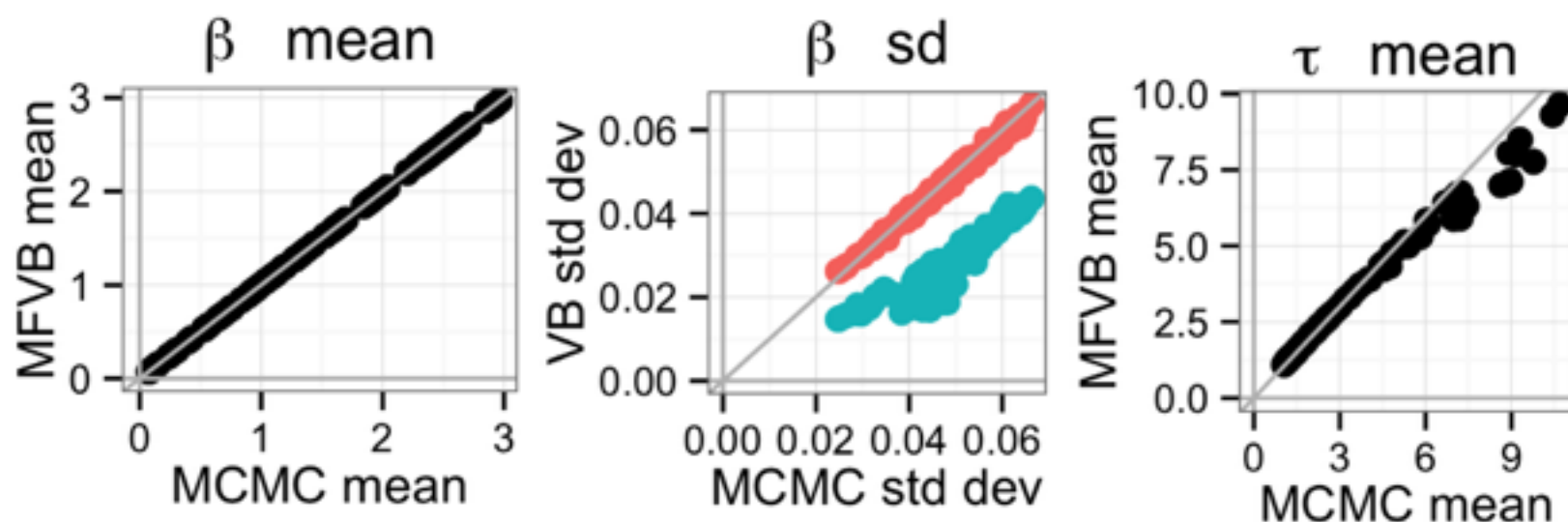
**LRVB**, **MFVB**



# Experiments

- Non-conjugate normal-Poisson generalized linear mixed model  
$$z_n | \beta, \tau \stackrel{\text{indep}}{\sim} \mathcal{N}(z_n | \beta x_n, \tau^{-1}), \quad y_n | z_n \stackrel{\text{indep}}{\sim} \text{Poisson}(y_n | \exp(z_n)),$$
$$\beta \sim \mathcal{N}(\beta | 0, \sigma_\beta^2), \quad \tau \sim \text{Gamma}(\tau | \alpha_\tau, \beta_\tau)$$
- 100 simulated data sets, 500 data points each, R MCMCglmm package

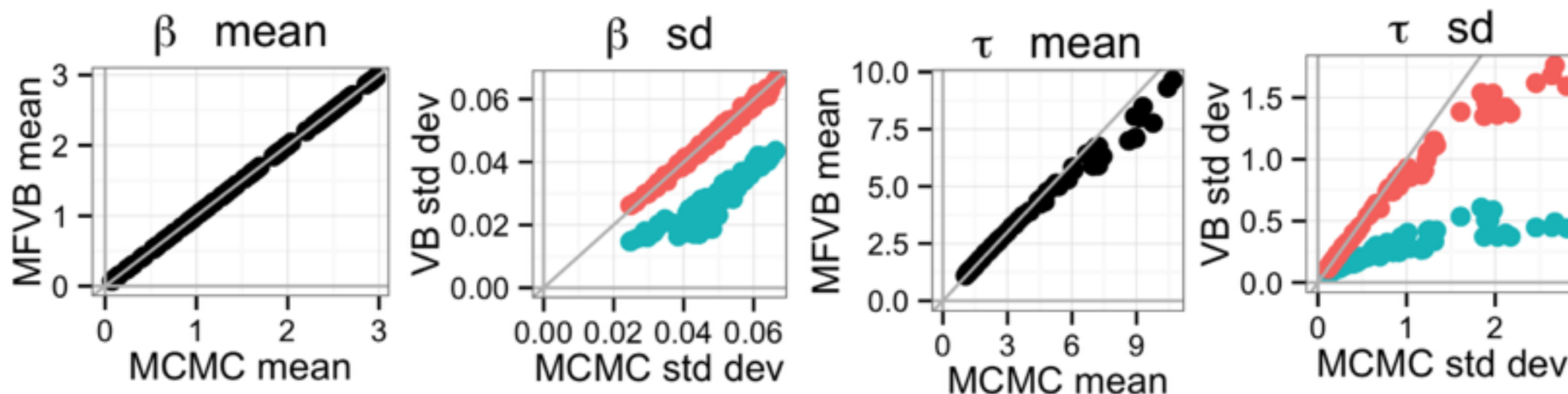
**LRVB**, **MFVB**



# Experiments

- Non-conjugate normal-Poisson generalized linear mixed model  
$$z_n | \beta, \tau \stackrel{indep}{\sim} \mathcal{N}(z_n | \beta x_n, \tau^{-1}), \quad y_n | z_n \stackrel{indep}{\sim} \text{Poisson}(y_n | \exp(z_n)),$$
$$\beta \sim \mathcal{N}(\beta | 0, \sigma_\beta^2), \quad \tau \sim \text{Gamma}(\tau | \alpha_\tau, \beta_\tau)$$
- 100 simulated data sets, 500 data points each, R MCMCglmm package

**LRVB**, **MFVB**



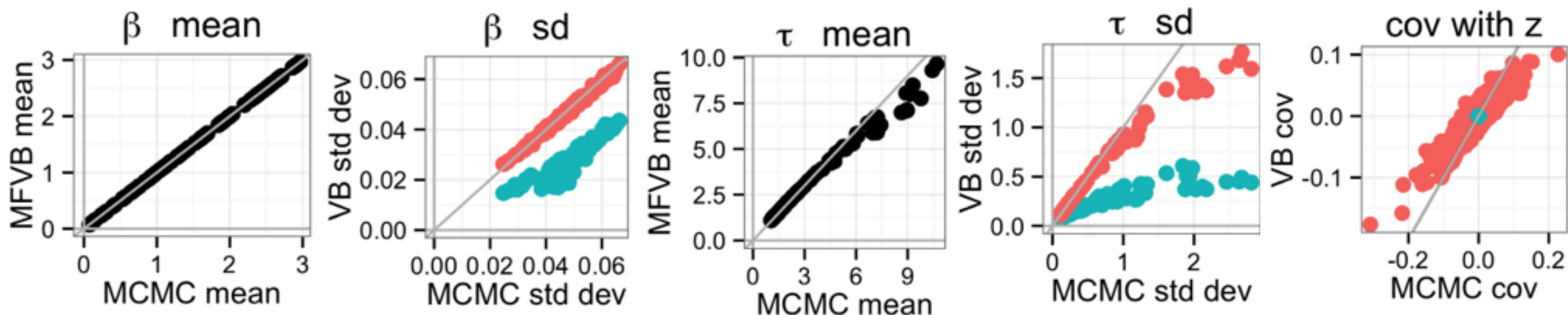
# Experiments

- Non-conjugate normal-Poisson generalized linear mixed model  

$$z_n | \beta, \tau \stackrel{indep}{\sim} \mathcal{N}(z_n | \beta x_n, \tau^{-1}), \quad y_n | z_n \stackrel{indep}{\sim} \text{Poisson}(y_n | \exp(z_n)),$$

$$\beta \sim \mathcal{N}(\beta | 0, \sigma_\beta^2), \quad \tau \sim \text{Gamma}(\tau | \alpha_\tau, \beta_\tau)$$
- 100 simulated data sets, 500 data points each, R MCMCglmm package

**LRVB**, **MFVB**



# Scaling the matrix inverse

# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$

# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$
- Decomposition of parameter vector

# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$
- Decomposition of parameter vector

$$\theta = (\alpha^T, z^T)^T$$



# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$
- Decomposition of parameter vector

$$\theta = (\alpha^T, z^T)^T$$

$$H =$$

$H_{\alpha}$	$H_{\alpha z}$
$H_{z\alpha}$	$H_z$

# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$

- Decomposition of parameter vector

$$\theta = (\alpha^T, z^T)^T$$

$$H = \begin{array}{|c|c|} \hline H_{\alpha} & H_{\alpha z} \\ \hline H_{z\alpha} & H_z \\ \hline \end{array}$$

- Schur complement

# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$

- Decomposition of parameter vector

$$\theta = (\alpha^T, z^T)^T$$

$$H = \begin{array}{|c|c|} \hline H_{\alpha} & H_{\alpha z} \\ \hline H_{z\alpha} & H_z \\ \hline \end{array}$$

- Schur complement

$$\hat{\Sigma}_{\alpha} = (I_{\alpha} - V_{\alpha}H_{\alpha} - V_{\alpha}H_{\alpha z} (I_z - V_zH_z)^{-1}V_zH_{z\alpha})^{-1} V_{\alpha}$$

# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$

- Decomposition of parameter vector

$$\theta = (\alpha^T, z^T)^T$$

$$H = \begin{array}{|c|c|} \hline H_{\alpha} & H_{\alpha z} \\ \hline H_{z\alpha} & H_z \\ \hline \end{array}$$

- Schur complement

$$\hat{\Sigma}_{\alpha} = (I_{\alpha} - V_{\alpha}H_{\alpha} - V_{\alpha}H_{\alpha z} (I_z - V_zH_z)^{-1}V_zH_{z\alpha})^{-1}V_{\alpha}$$

# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$

- Decomposition of parameter vector

$$\theta = (\alpha^T, z^T)^T$$

$$H = \begin{array}{|c|c|} \hline H_{\alpha} & H_{\alpha z} \\ \hline H_{z\alpha} & H_z \\ \hline \end{array}$$

- Schur complement

$$\hat{\Sigma}_{\alpha} = (I_{\alpha} - V_{\alpha}H_{\alpha} - V_{\alpha}H_{\alpha z} (I_z - V_z H_z)^{-1} V_z H_{z\alpha})^{-1} V_{\alpha}$$

# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$

- Decomposition of parameter vector

$$\theta = (\alpha^T, z^T)^T$$

$$H = \begin{array}{|c|c|} \hline H_{\alpha} & H_{\alpha z} \\ \hline H_{z\alpha} & H_z \\ \hline \end{array}$$

- Schur complement

$$\hat{\Sigma}_{\alpha} = (I_{\alpha} - V_{\alpha}H_{\alpha} - V_{\alpha}H_{\alpha z} (I_z - V_zH_z)^{-1}V_zH_{z\alpha})^{-1} V_{\alpha}$$

- Sparsity patterns

# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$

- Decomposition of parameter vector

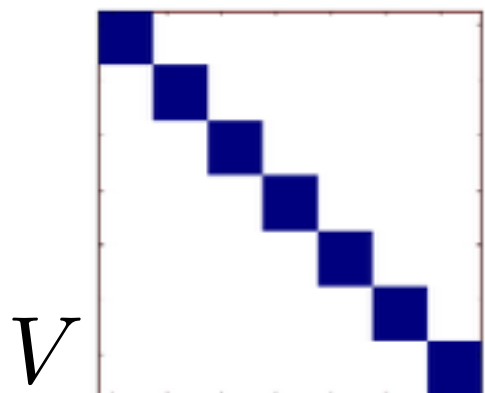
$$\theta = (\alpha^T, z^T)^T$$

$$H = \begin{bmatrix} H_{\alpha} & H_{\alpha z} \\ H_{z\alpha} & H_z \end{bmatrix}$$

- Schur complement

$$\hat{\Sigma}_{\alpha} = (I_{\alpha} - V_{\alpha}H_{\alpha} - V_{\alpha}H_{\alpha z} (I_z - V_zH_z)^{-1}V_zH_{z\alpha})^{-1} V_{\alpha}$$

- Sparsity patterns



# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$

- Decomposition of parameter vector

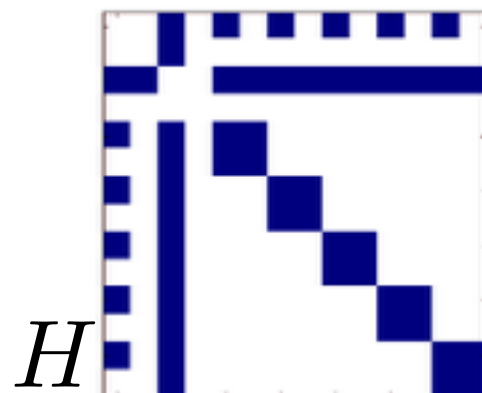
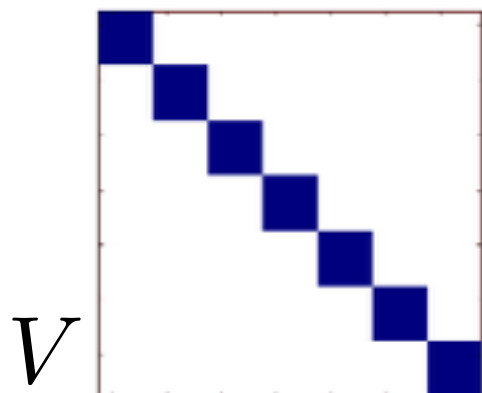
$$\theta = (\alpha^T, z^T)^T$$

$$H = \begin{bmatrix} H_{\alpha} & H_{\alpha z} \\ H_{z\alpha} & H_z \end{bmatrix}$$

- Schur complement

$$\hat{\Sigma}_{\alpha} = (I_{\alpha} - V_{\alpha}H_{\alpha} - V_{\alpha}H_{\alpha z} (I_z - V_zH_z)^{-1}V_zH_{z\alpha})^{-1} V_{\alpha}$$

- Sparsity patterns





# Scaling the matrix inverse

- LRVB estimate  $\hat{\Sigma} = (I - VH)^{-1}V$

- Decomposition of parameter vector

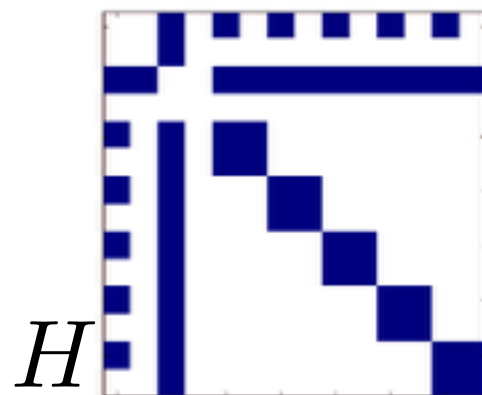
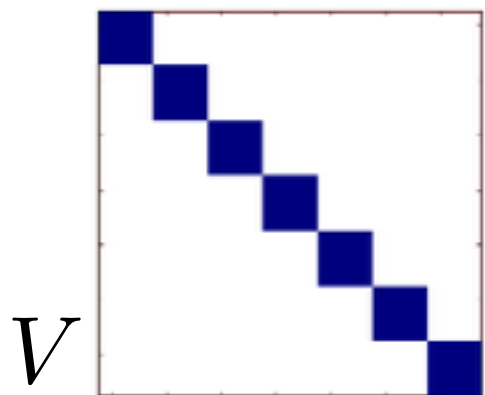
$$\theta = (\alpha^T, z^T)^T$$

$$H = \begin{bmatrix} H_{\alpha} & H_{\alpha z} \\ H_{z\alpha} & H_z \end{bmatrix}$$

- Schur complement

$$\hat{\Sigma}_{\alpha} = (I_{\alpha} - V_{\alpha}H_{\alpha} - V_{\alpha}H_{\alpha z} (I_z - V_zH_z)^{-1}V_zH_{z\alpha})^{-1} V_{\alpha}$$

- Sparsity patterns



# Roadmap

- Variational Bayes as an alternative to MCMC
- Challenges of VB
- Accurate uncertainties from VB
- Accurate robustness quantification from VB

# Roadmap

- Variational Bayes as an alternative to MCMC
- Challenges of VB
- Accurate uncertainties from VB
- Accurate robustness quantification from VB

# Robustness quantification

- Bayes Theorem

$$p(\theta|x)$$

$$\propto_{\theta} p(x|\theta)p(\theta)$$

# Robustness quantification

- Bayes Theorem

$$p(\theta|x, \alpha)$$

$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

# Robustness quantification

- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

# Robustness quantification

- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

- Sensitivity

# Robustness quantification

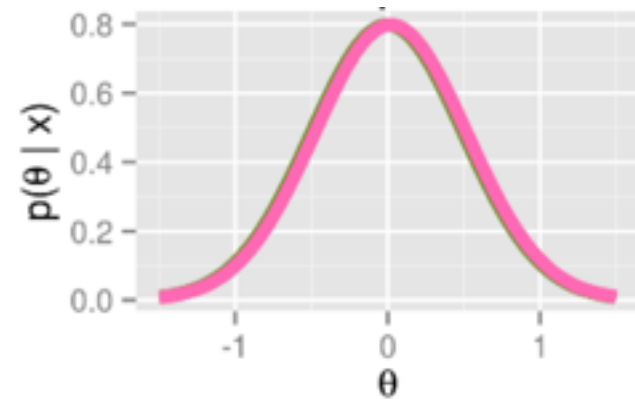
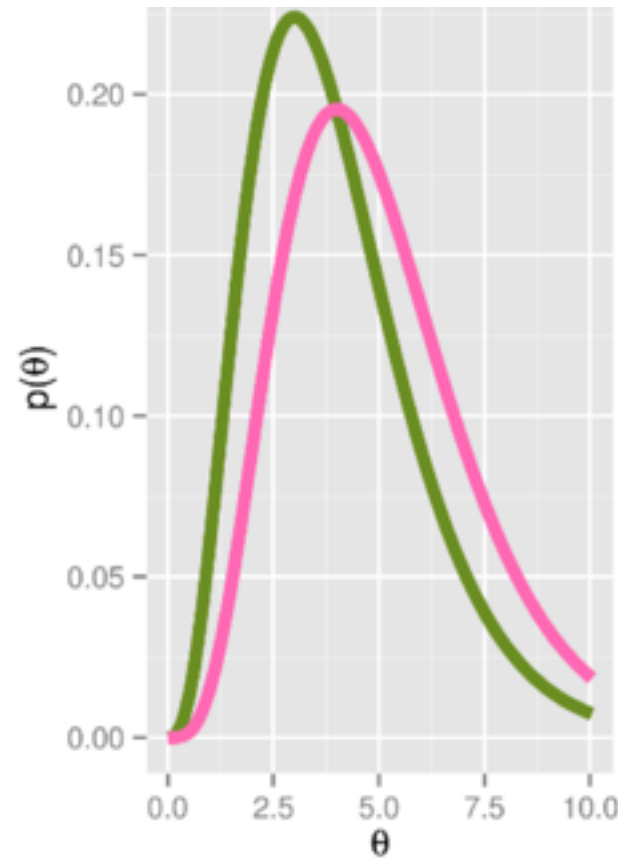
- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$

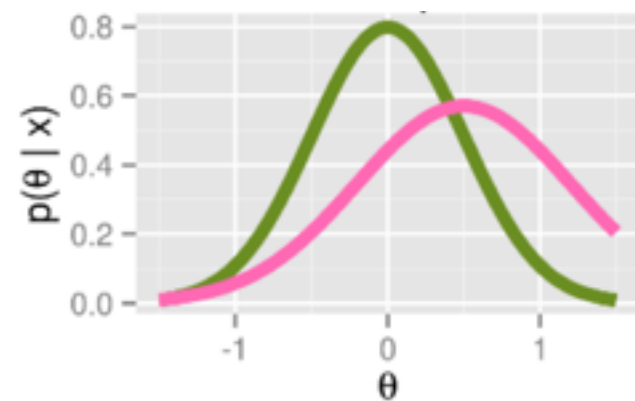
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

- Sensitivity

Some reasonable priors



Bayes Theorem





# Robustness quantification

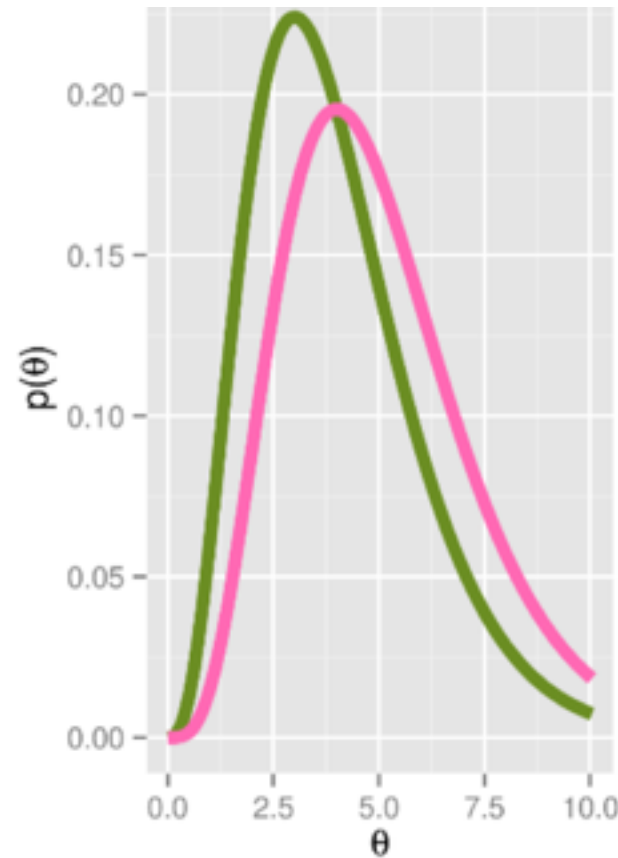
- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

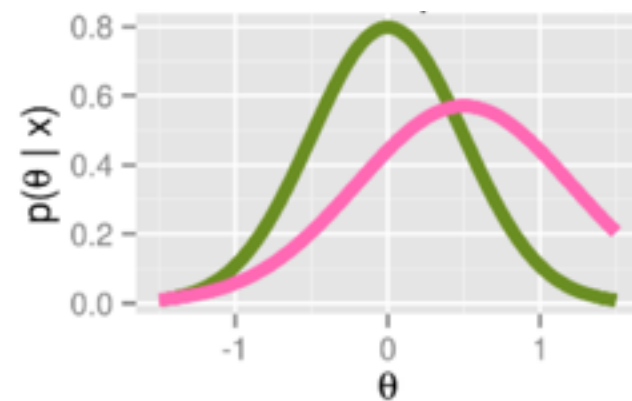
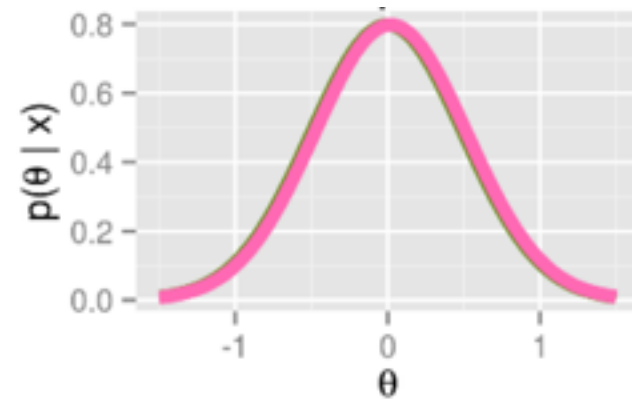
- Sensitivity

$$\mathbb{E}_{p_{\alpha}}[g(\theta)]$$

Some reasonable priors



Bayes Theorem



# Robustness quantification

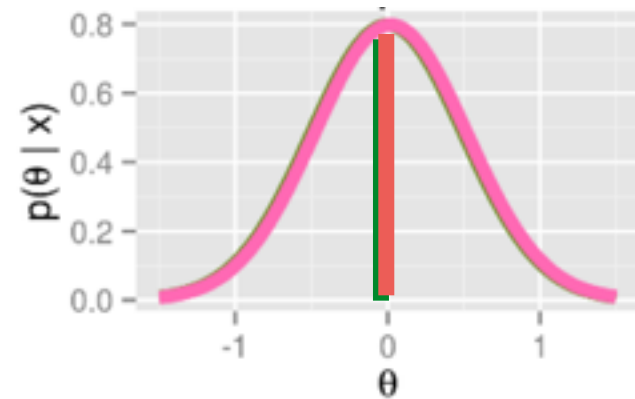
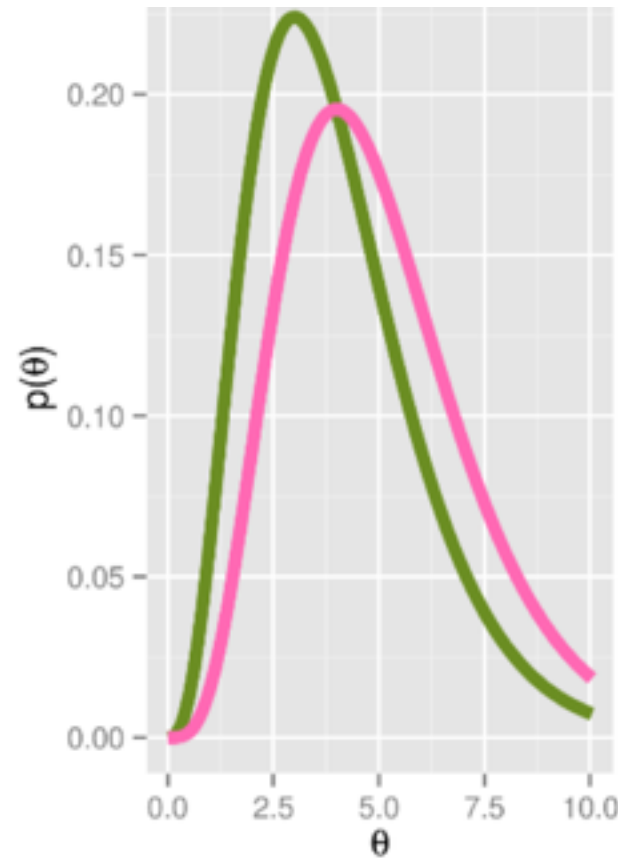
- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

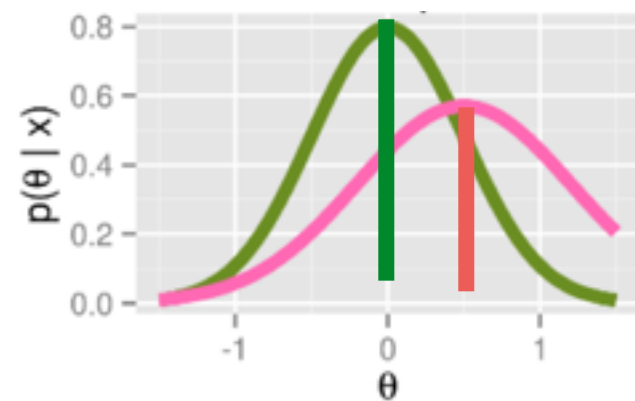
- Sensitivity

$$\mathbb{E}_{p_{\alpha}}[g(\theta)]$$

Some reasonable priors



Bayes Theorem



# Robustness quantification

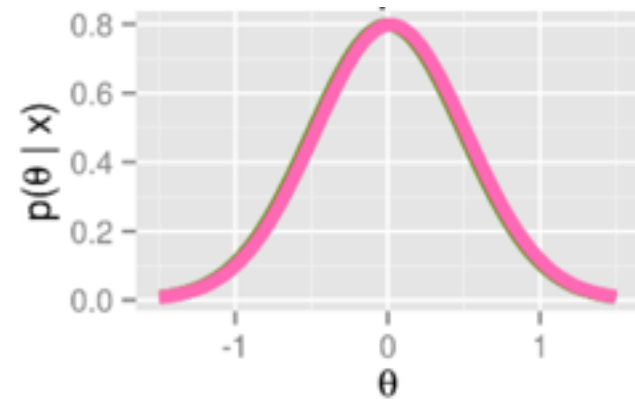
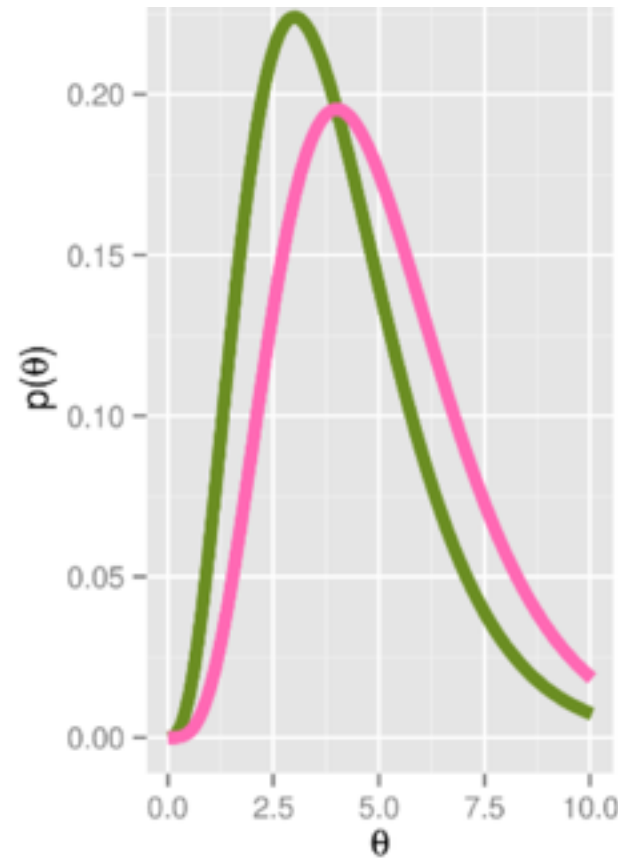
- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

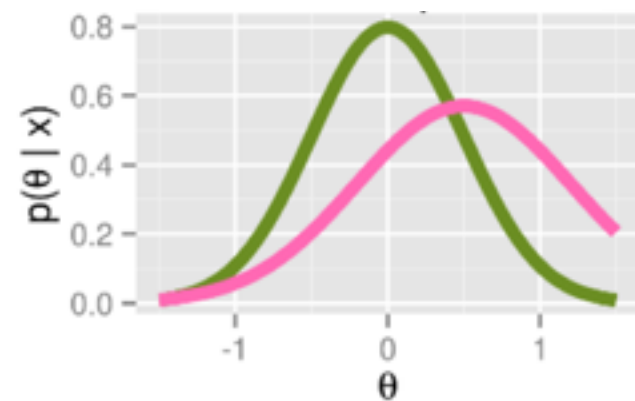
- Sensitivity

$$\mathbb{E}_{p_{\alpha}}[g(\theta)]$$

Some reasonable priors



Bayes Theorem



# Robustness quantification

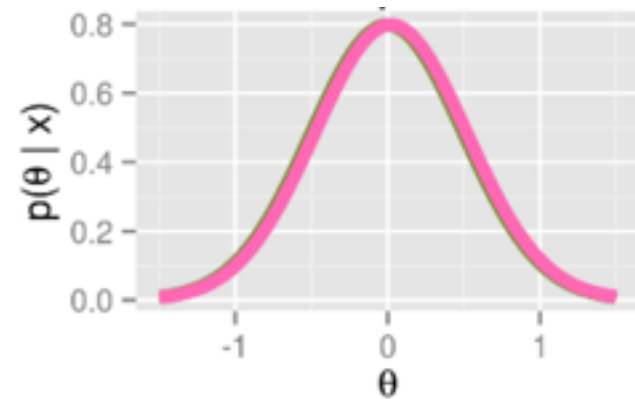
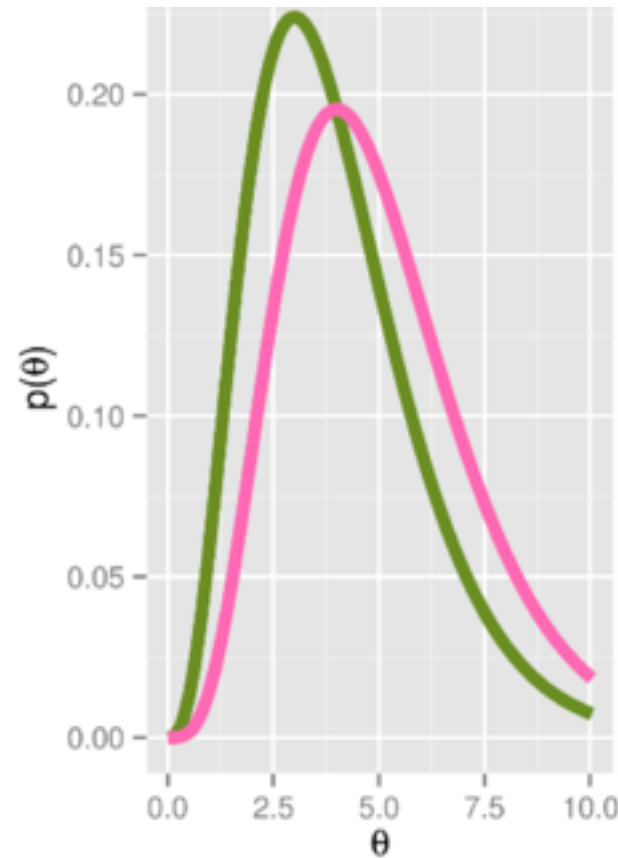
- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

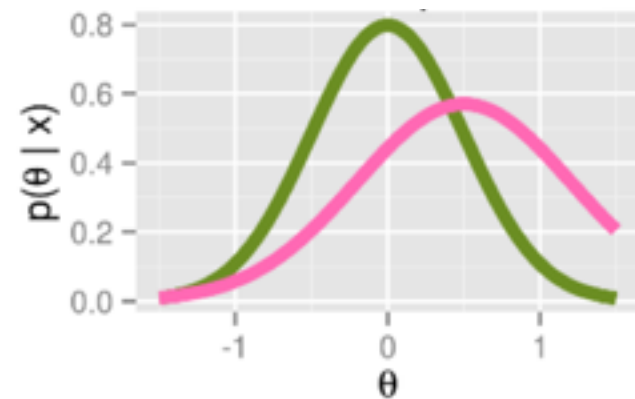
- Sensitivity

$$S := \left. \frac{d\mathbb{E}_{p_{\alpha}}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha$$

Some reasonable priors



Bayes Theorem



# Robustness quantification

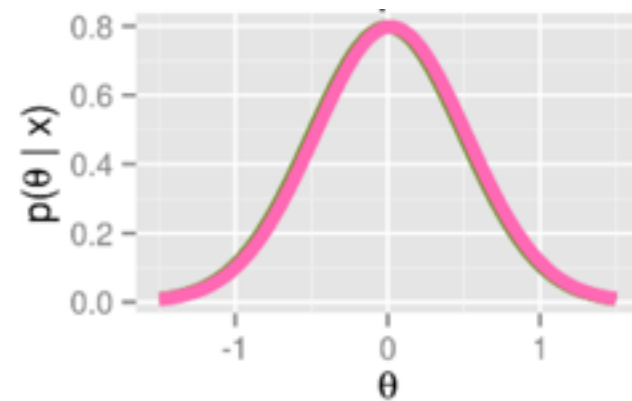
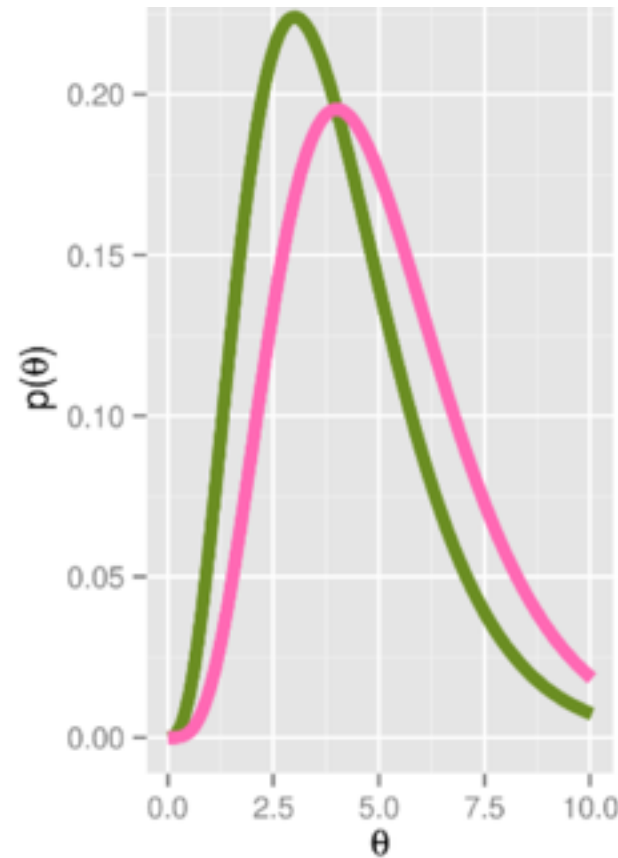
- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

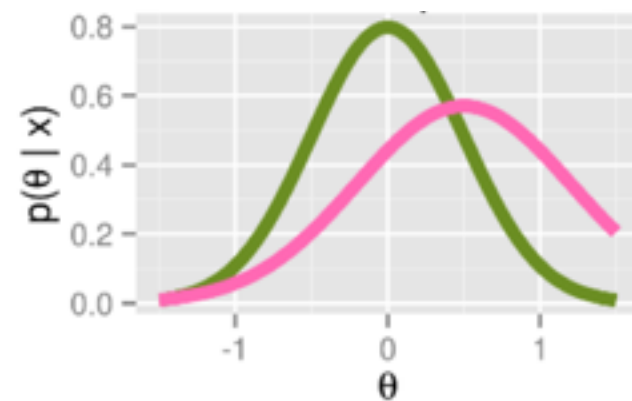
- Sensitivity

$$S := \left. \frac{d\mathbb{E}_{p_{\alpha}}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha$$

Some reasonable priors



Bayes Theorem



# Robustness quantification

- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$

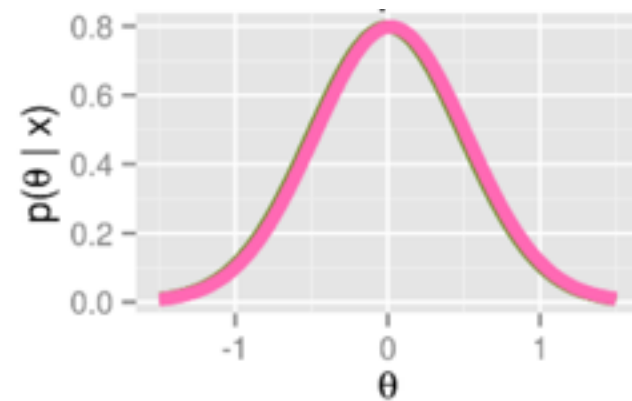
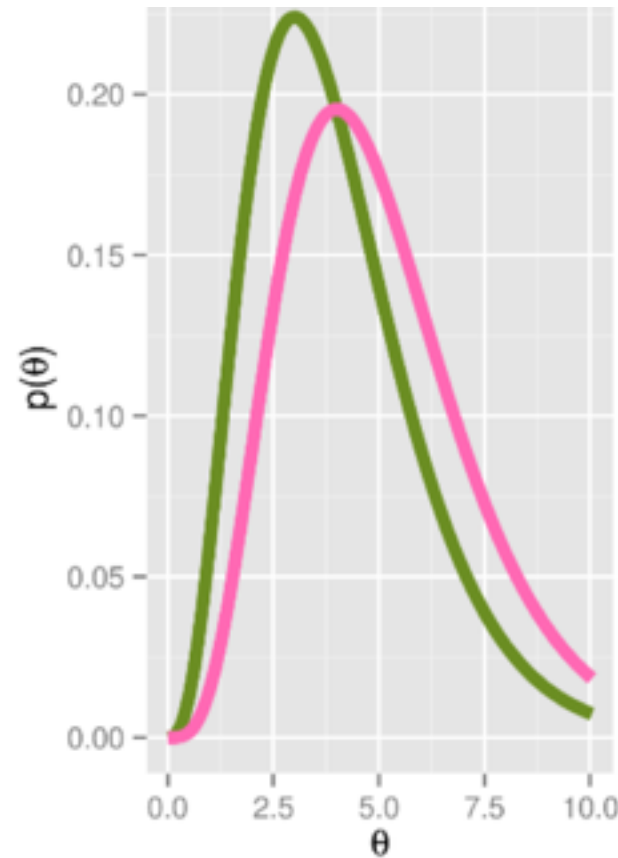
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

- Sensitivity

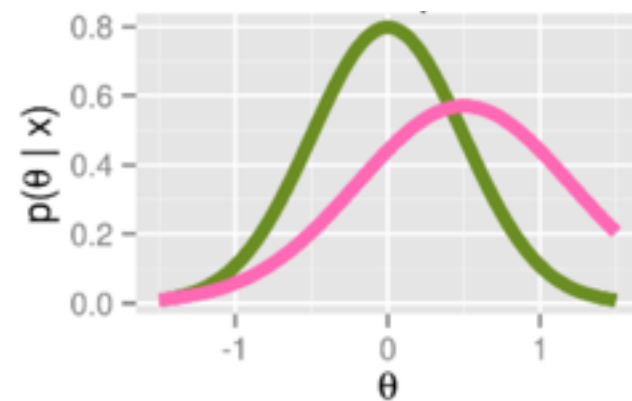
$$S := \left. \frac{d\mathbb{E}_{p_{\alpha}}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha$$

$$\approx \left. \frac{d\mathbb{E}_{q_{\alpha}^*}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha =: \hat{S}$$

Some reasonable priors



Bayes Theorem



# Robustness quantification

- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$

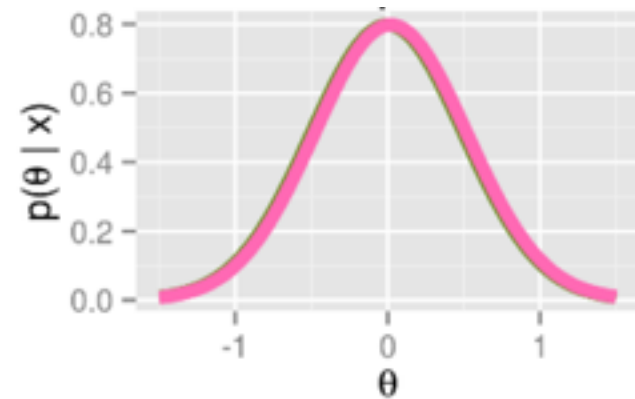
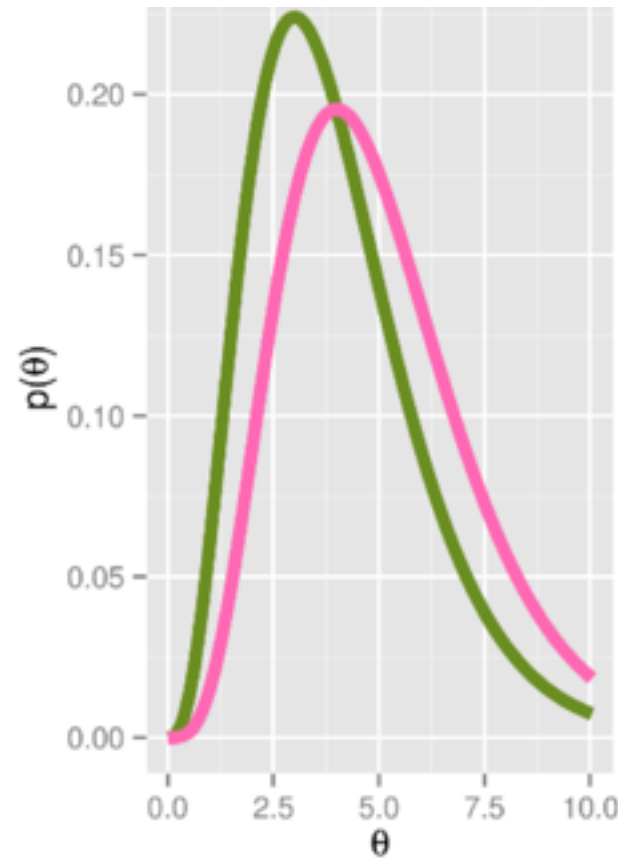
$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

- Sensitivity

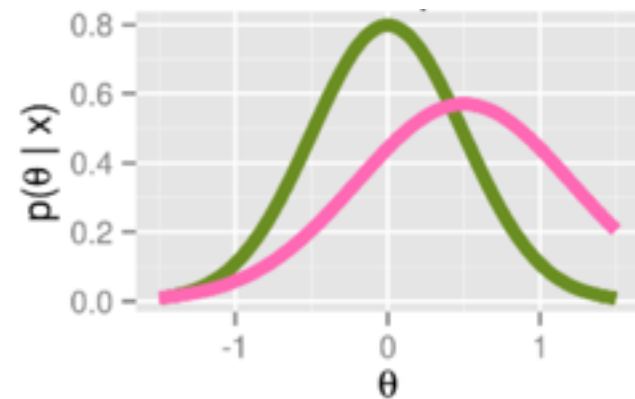
$$S := \left. \frac{d\mathbb{E}_{p_{\alpha}}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha$$

$$\approx \left. \frac{d\mathbb{E}_{q_{\alpha}^*}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha =: \hat{S}$$

Some reasonable priors



Bayes Theorem



LRVB estimator

# Robustness quantification

- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$

$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

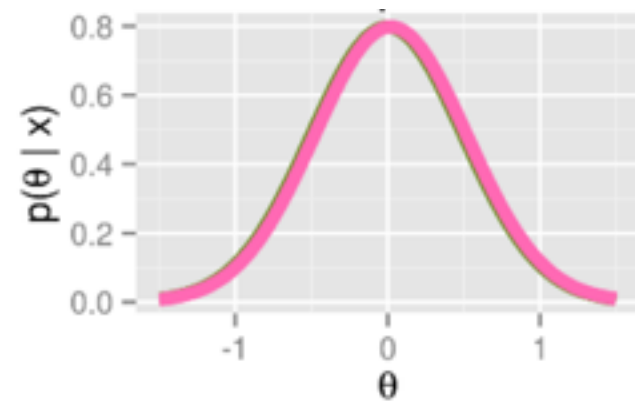
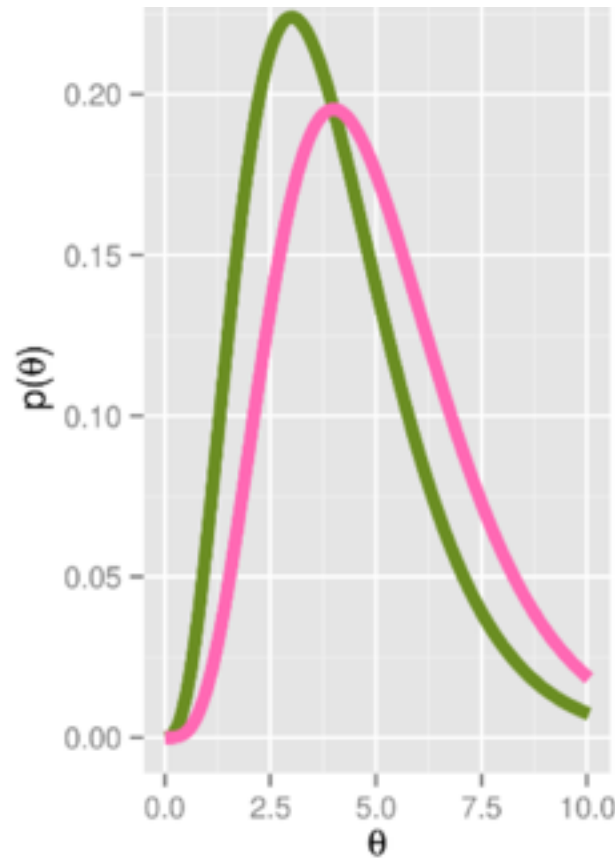
- Sensitivity

$$S := \left. \frac{d\mathbb{E}_{p_{\alpha}}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha$$

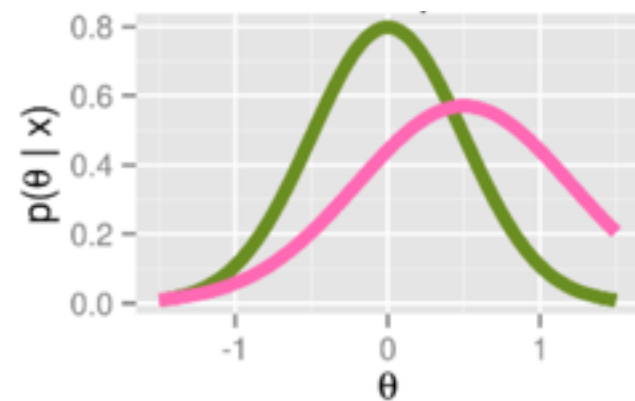
$$\approx \left. \frac{d\mathbb{E}_{q_{\alpha}^*}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha =: \hat{S}$$

- When  $q_{\alpha}^*$  in exponential family

Some reasonable priors



Bayes Theorem



LRVB estimator



# Robustness quantification

- Bayes Theorem

$$p_{\alpha}(\theta) := p(\theta|x, \alpha)$$

$$\propto_{\theta} p(x|\theta)p(\theta|\alpha)$$

- Sensitivity

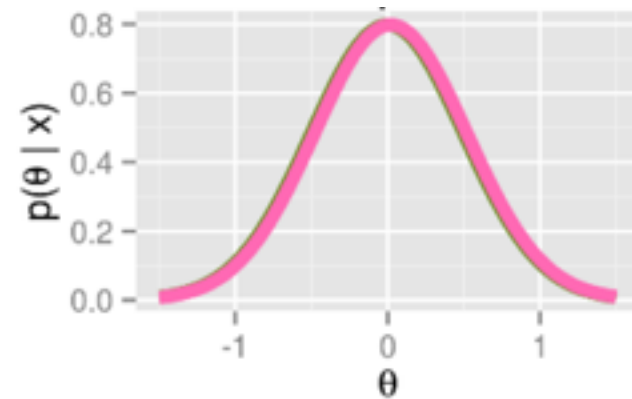
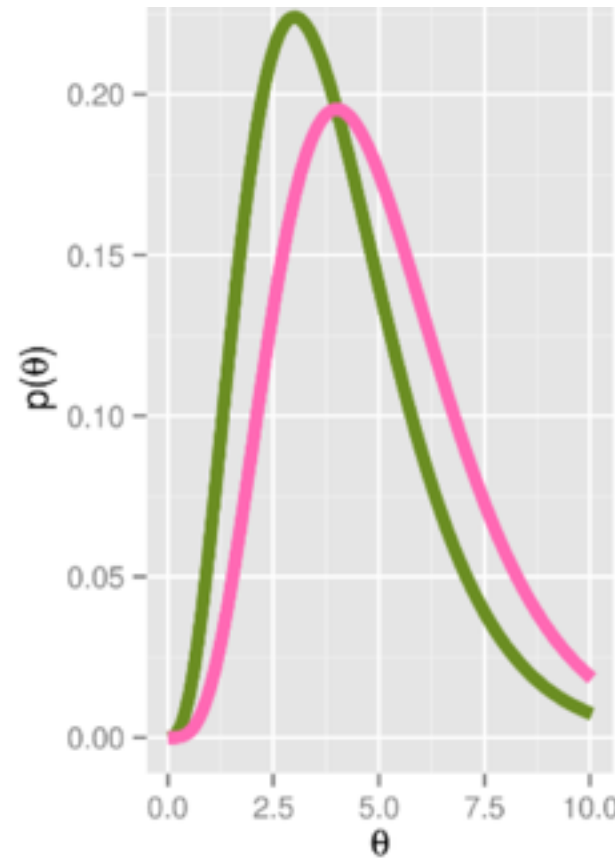
$$S := \left. \frac{d\mathbb{E}_{p_{\alpha}}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha$$

$$\approx \left. \frac{d\mathbb{E}_{q_{\alpha}^*}[g(\theta)]}{d\alpha} \right|_{\alpha} \Delta\alpha =: \hat{S} \quad \leftarrow \text{LRVB estimator}$$

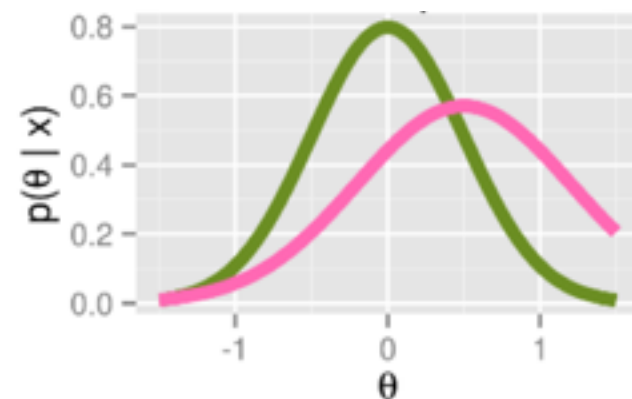
- When  $q_{\alpha}^*$  in exponential family

$$\hat{S} = A \left( \left. \frac{\partial^2 KL}{\partial m \partial m^T} \right|_{m=m^*} \right)^{-1} B$$

Some reasonable priors



Bayes Theorem



# Microcredit Experiment

- Simplified from Meager (2015)
- $K$  microcredit trials (Mexico, Mongolia, Bosnia, India, Morocco, Philippines, Ethiopia)
- $N_k$  businesses in  $k$ th site ( $\sim 900$  to  $\sim 17K$ )
- Profit of  $n$ th business at  $k$ th site:

profit  $\rightarrow y_{kn} \stackrel{indep}{\sim} \mathcal{N}(\mu_k + T_{kn}\tau_k, \sigma_k^2)$

1 if microcredit  $\rightarrow T_{kn}$

- Priors and hyperpriors:

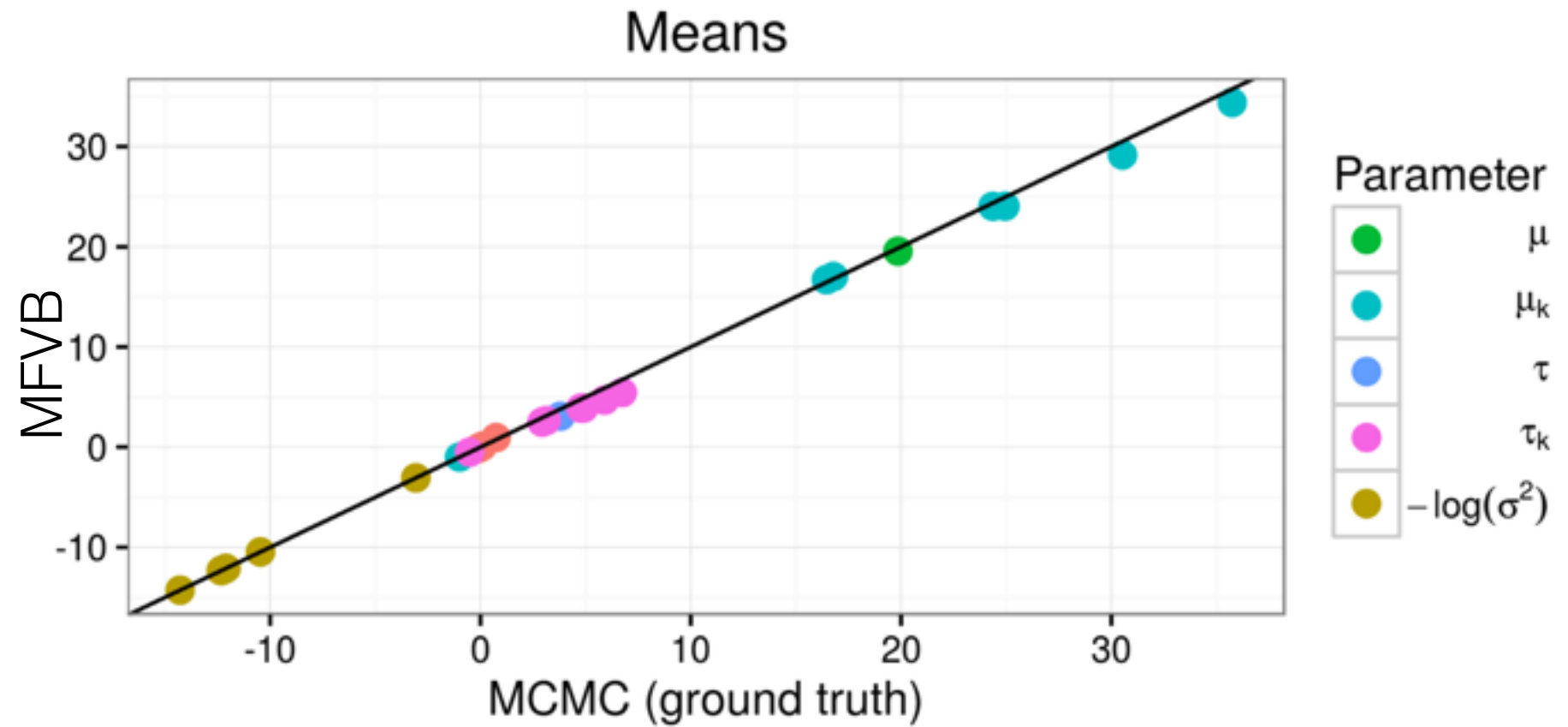
$$\begin{pmatrix} \mu_k \\ \tau_k \end{pmatrix} \stackrel{iid}{\sim} \mathcal{N}\left(\begin{pmatrix} \mu \\ \tau \end{pmatrix}, C\right) \quad \begin{pmatrix} \mu \\ \tau \end{pmatrix} \stackrel{iid}{\sim} \mathcal{N}\left(\begin{pmatrix} \mu_0 \\ \tau_0 \end{pmatrix}, \Lambda^{-1}\right)$$

$$\sigma_k^{-2} \stackrel{iid}{\sim} \Gamma(a, b)$$

$$C \sim \text{Sep\&LKJ}(\eta, c, d)$$

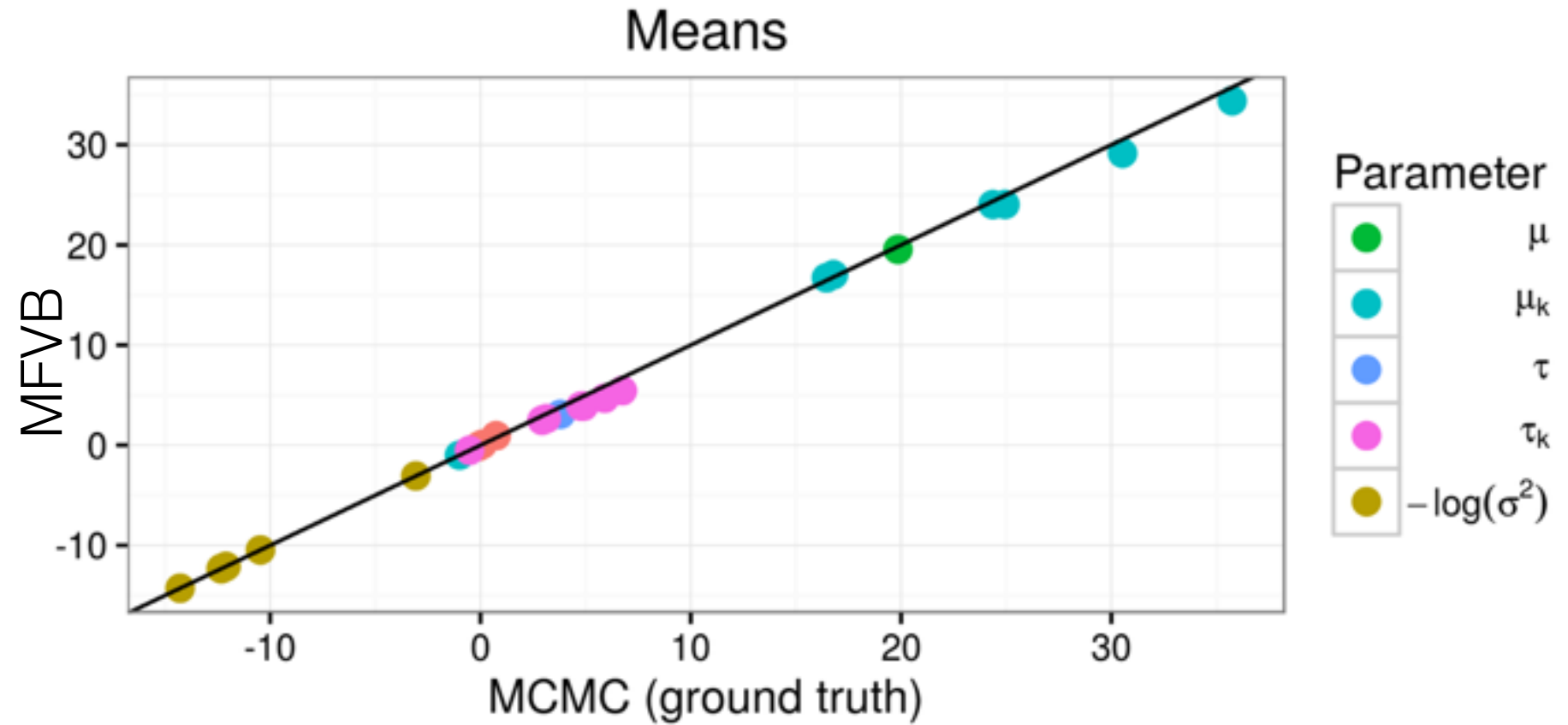
# Microcredit Experiment

# Microcredit Experiment



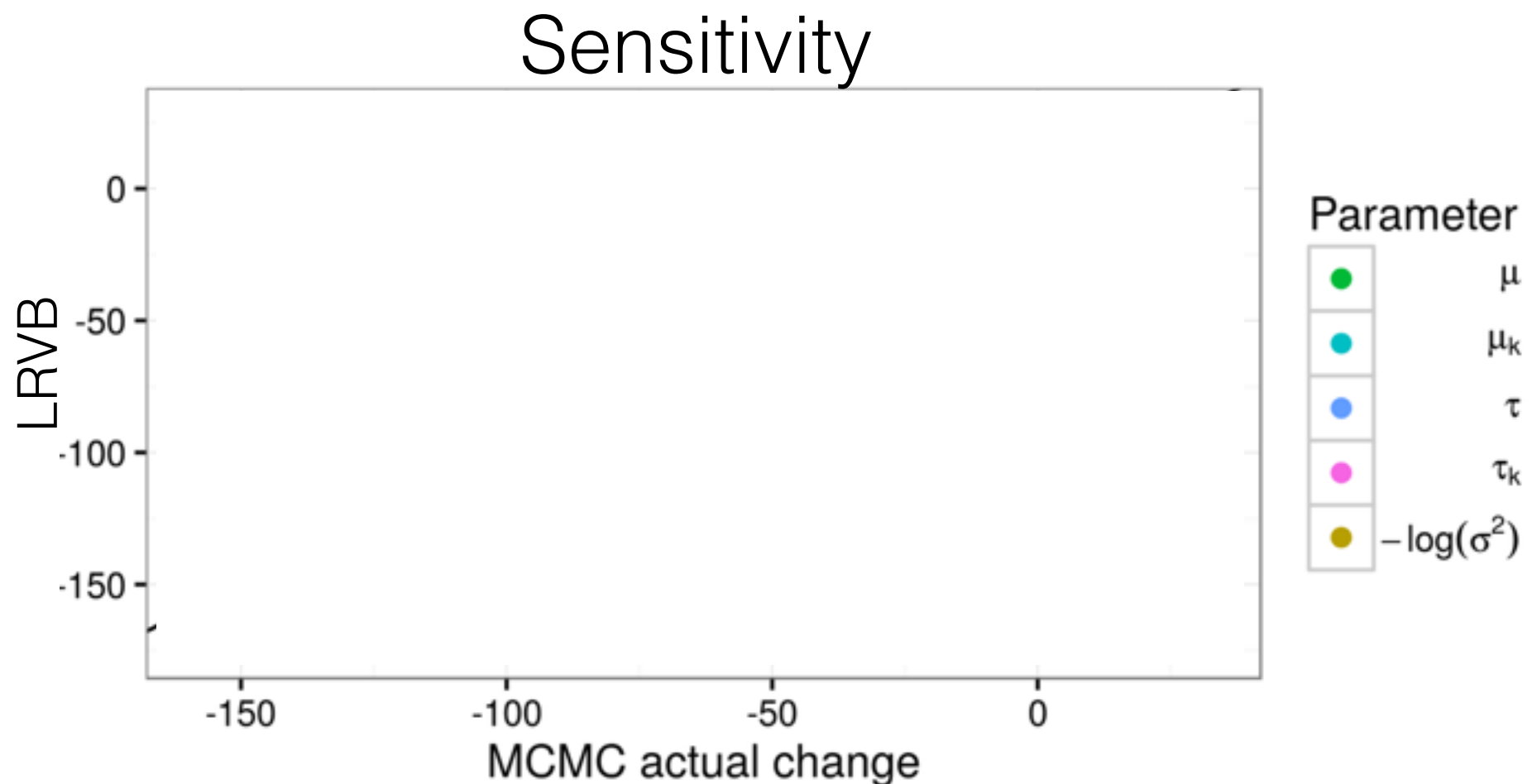
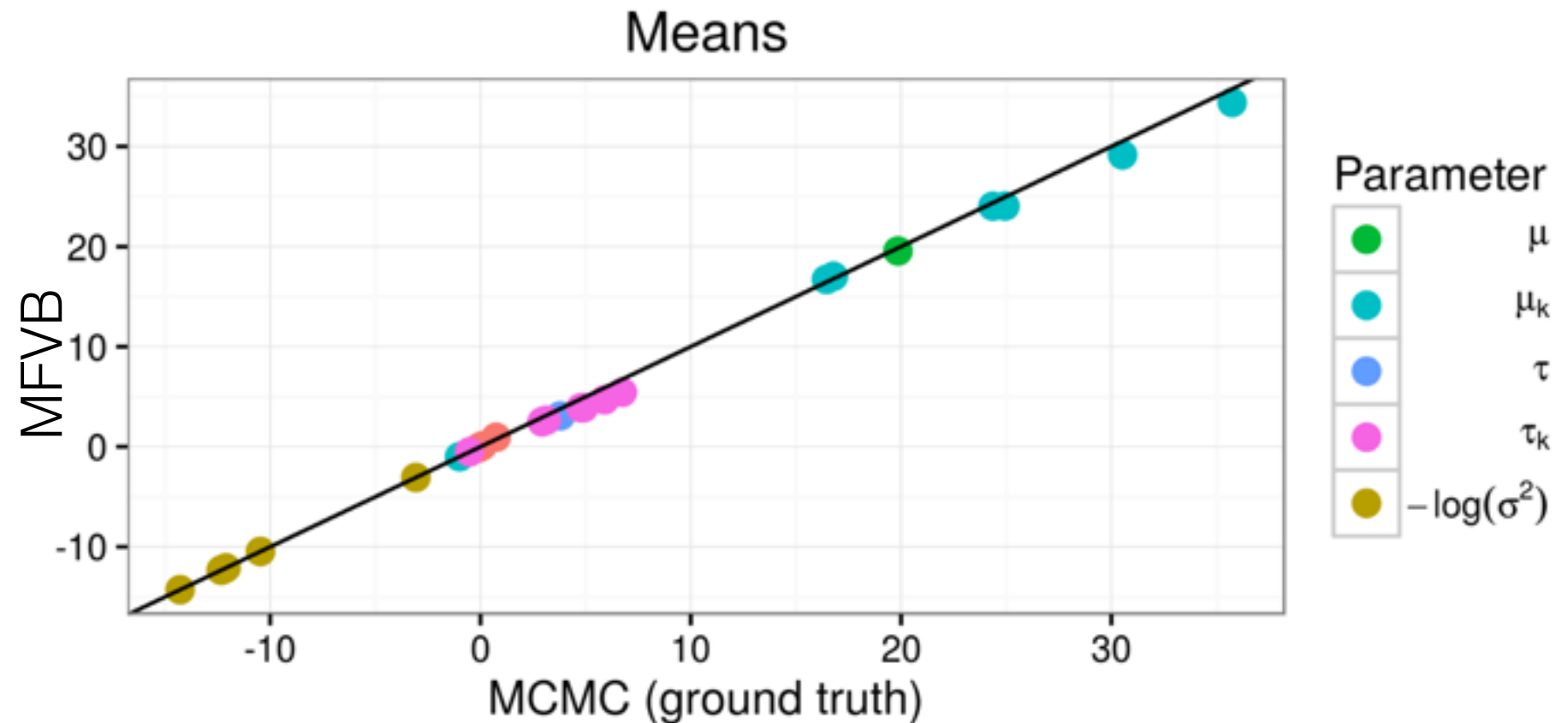
# Microcredit Experiment

- Perturb  $\Lambda_{11}$ :  
0.03  $\rightarrow$  0.04



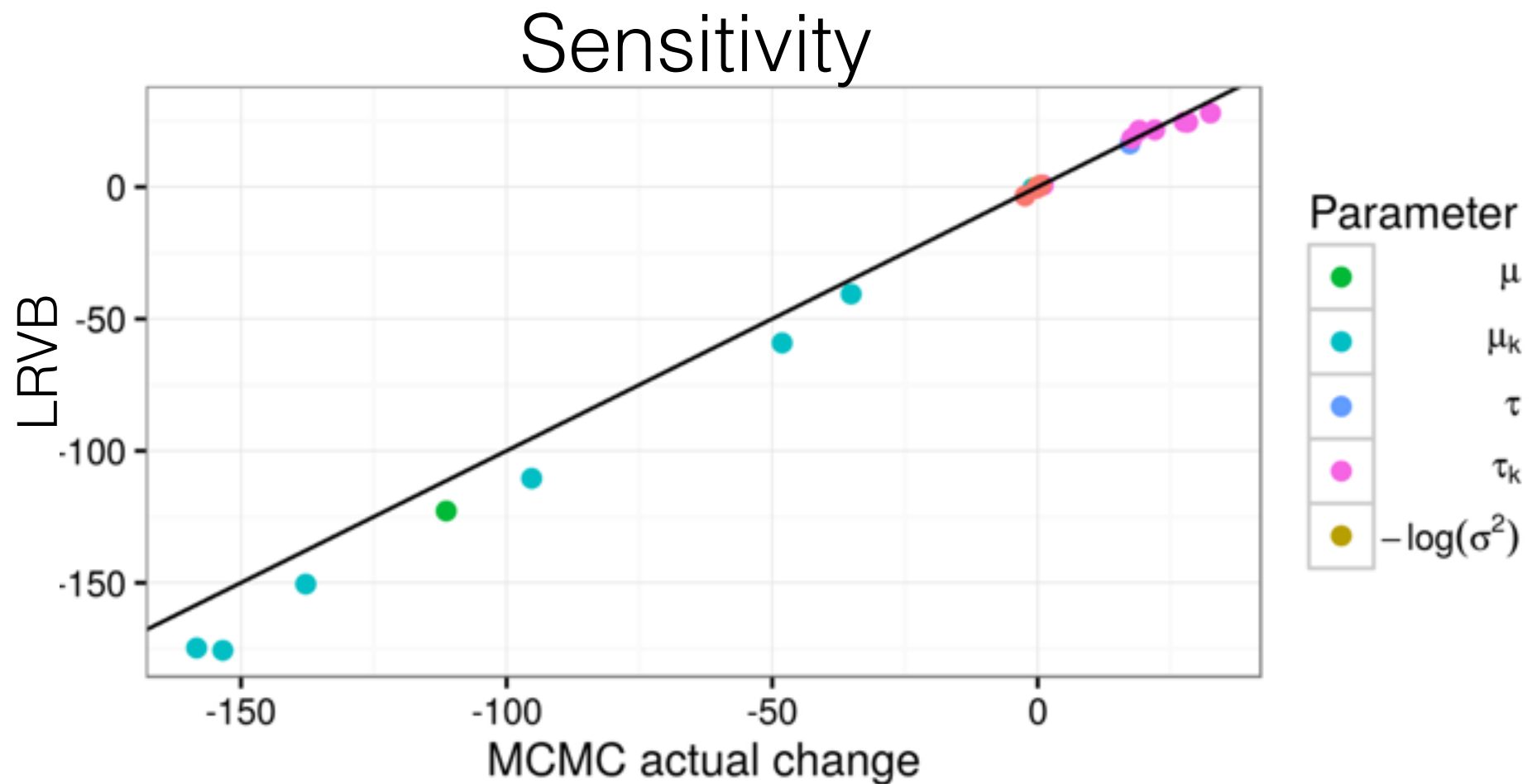
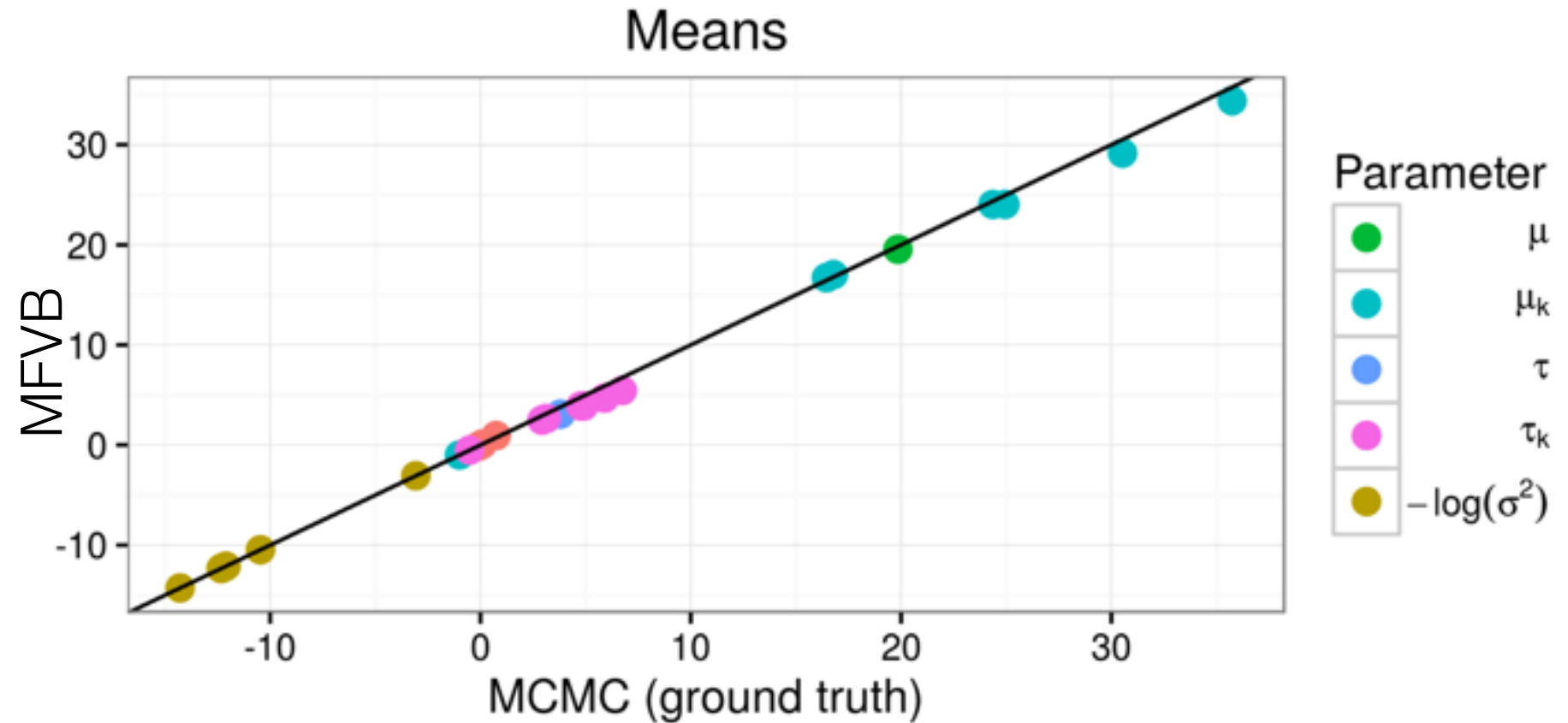
# Microcredit Experiment

- Perturb  $\Lambda_{11}$ :  
0.03  $\rightarrow$  0.04



# Microcredit Experiment

- Perturb  $\Lambda_{11}$ :  
 $0.03 \rightarrow 0.04$



# Microcredit Experiment



# Microcredit Experiment

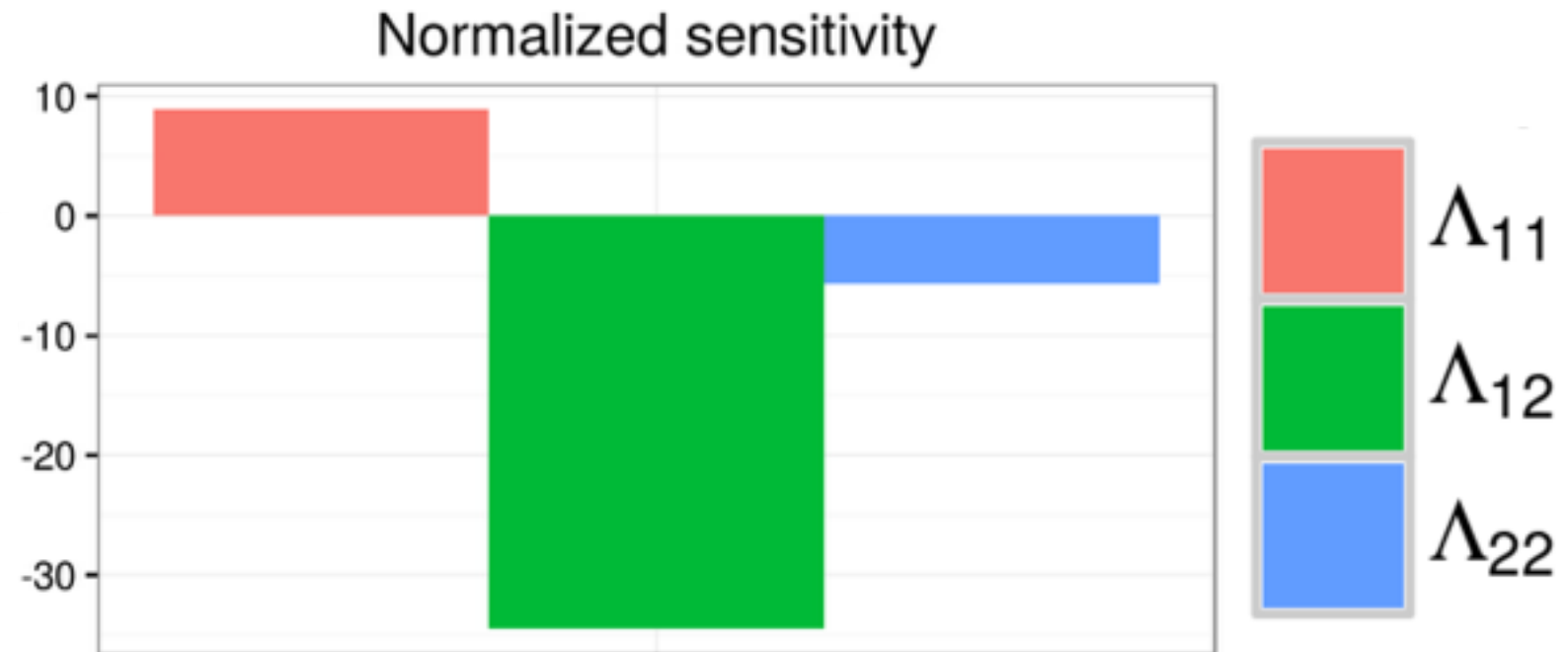
- Sensitivity of the expected microcredit effect ( $\tau$ )

# Microcredit Experiment

- Sensitivity of the expected microcredit effect ( $\tau$ )
- Normalized to be on scale of  $\tau$  std devs

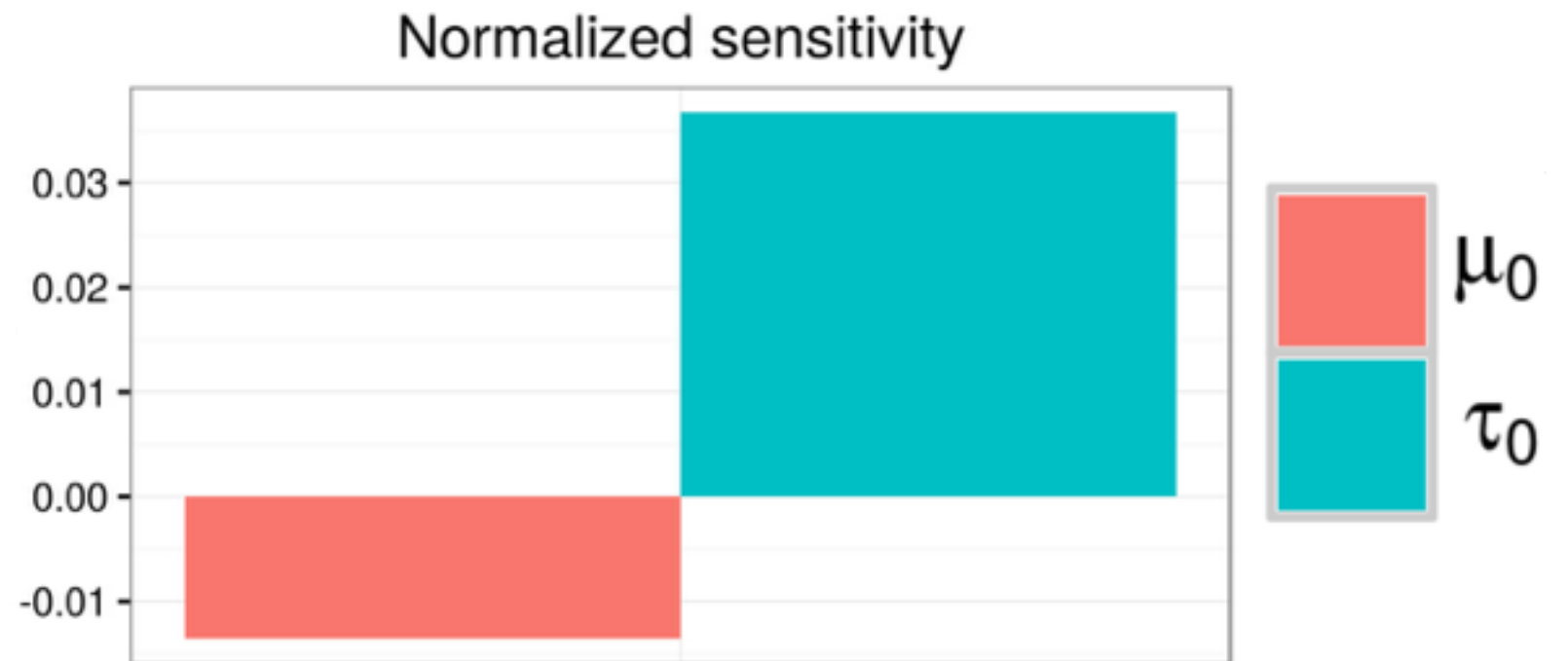
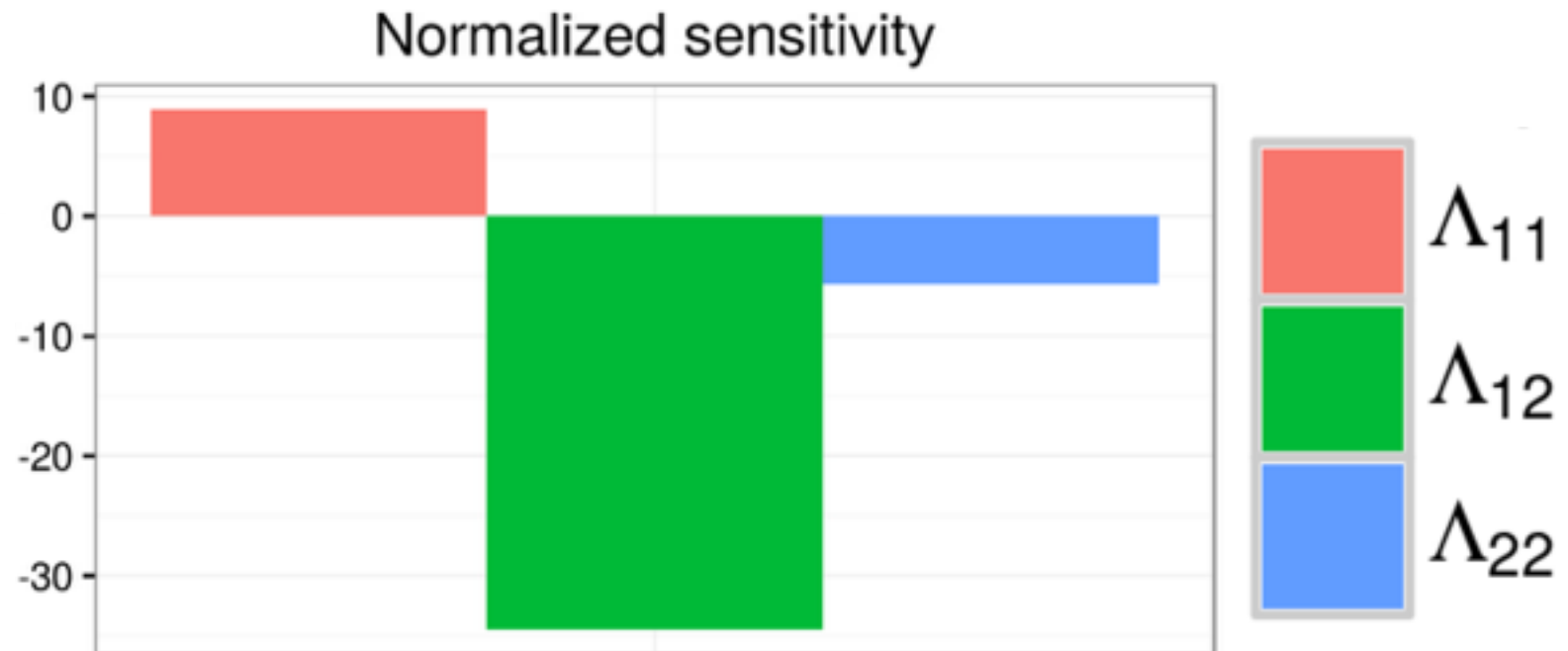
# Microcredit Experiment

- Sensitivity of the expected microcredit effect ( $\tau$ )
- Normalized to be on scale of  $\tau$  std devs



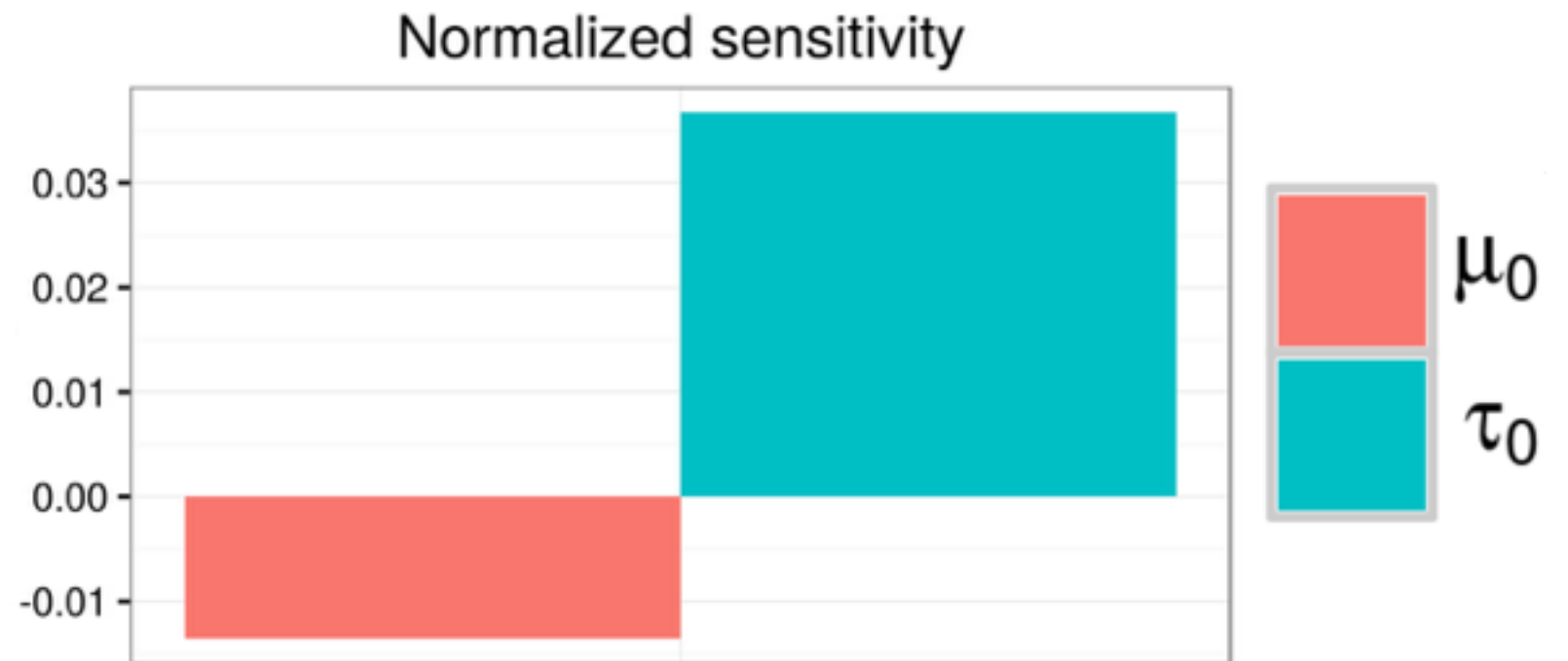
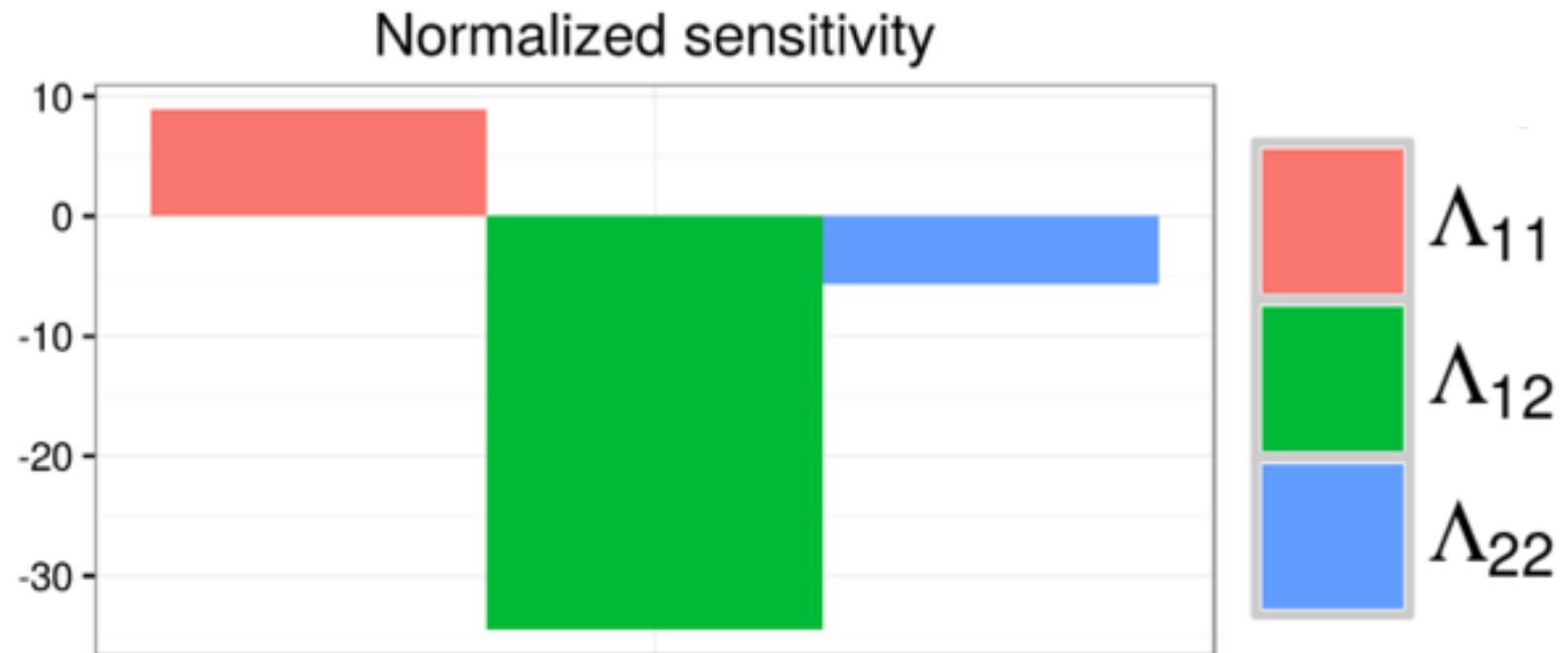
# Microcredit Experiment

- Sensitivity of the expected microcredit effect ( $\tau$ )
- Normalized to be on scale of  $\tau$  std devs



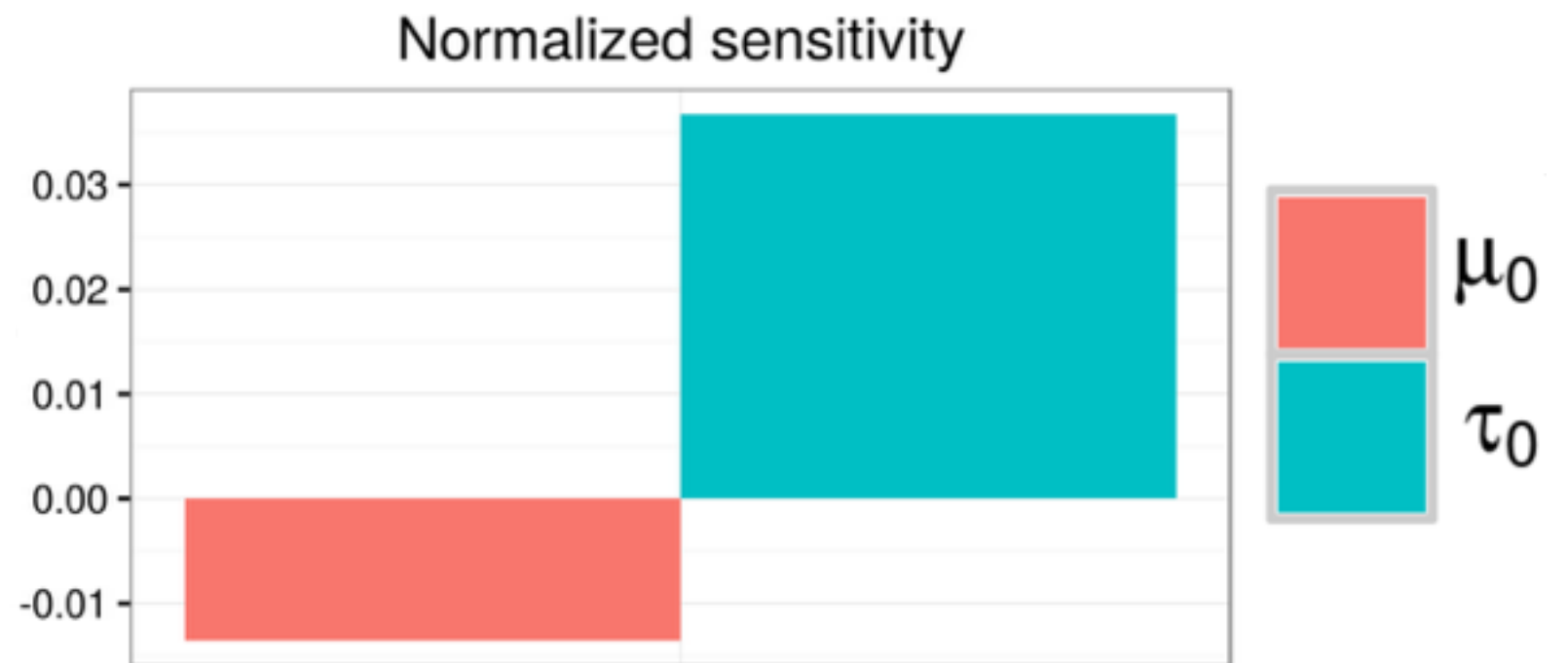
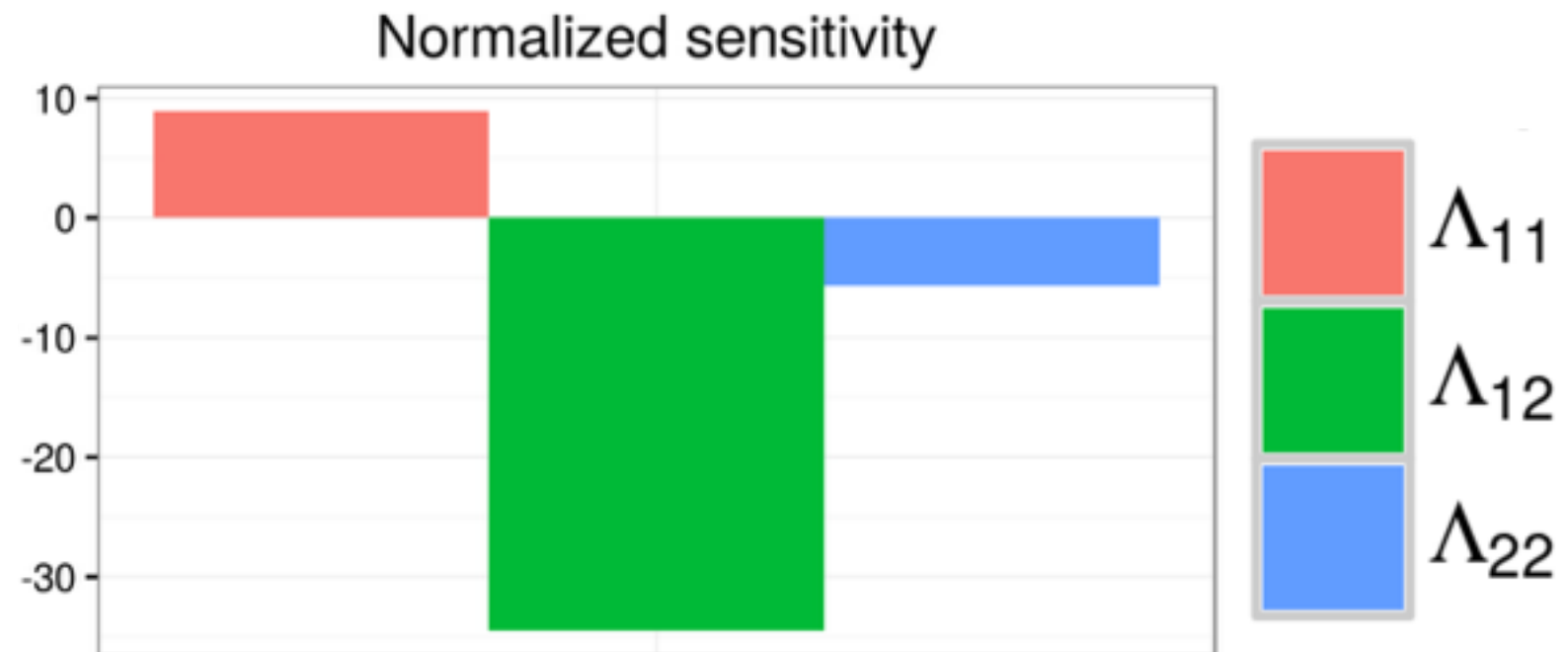
# Microcredit Experiment

- Sensitivity of the expected microcredit effect ( $\tau$ )
- Normalized to be on scale of  $\tau$  std devs
- $\tau$  mean (MFVB): 3.08 USD PPP
- $\tau$  std dev (LRVB): 1.83 USD PPP
- Mean is 1.68 std dev from 0



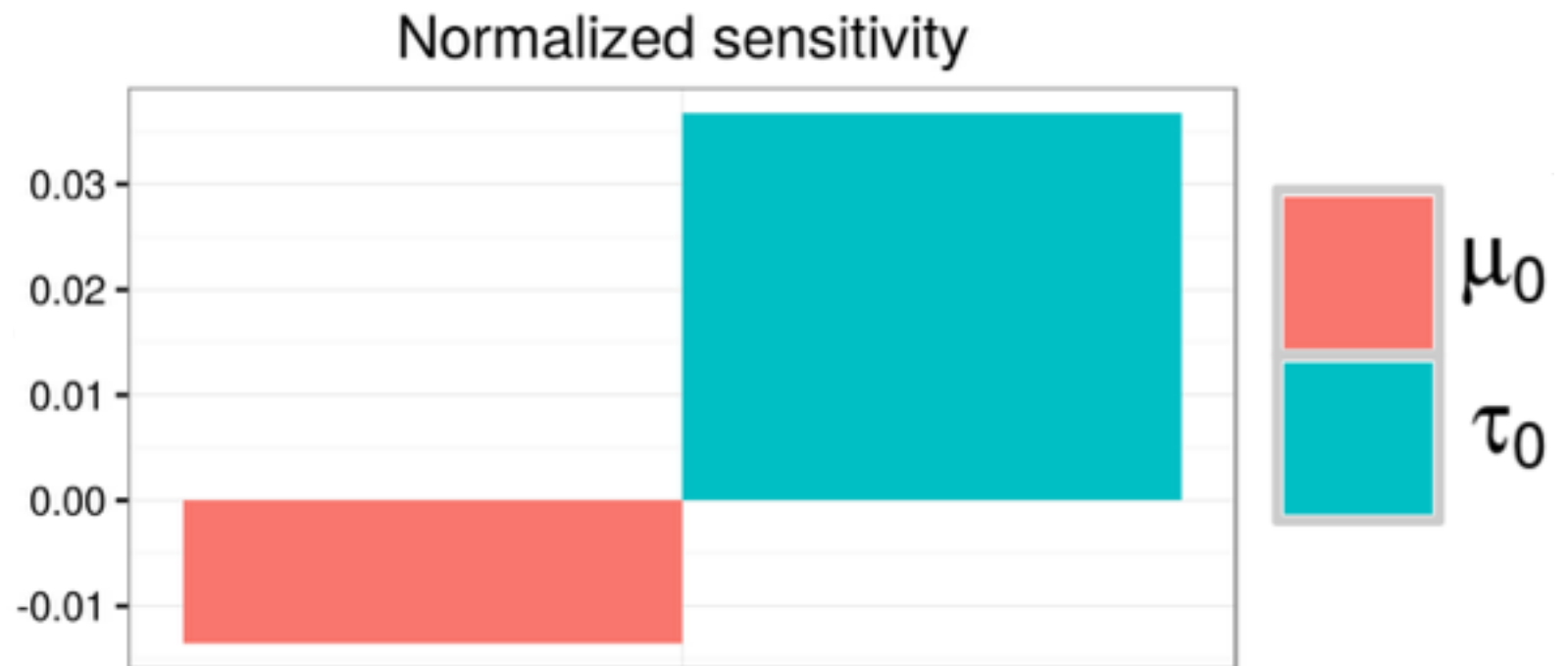
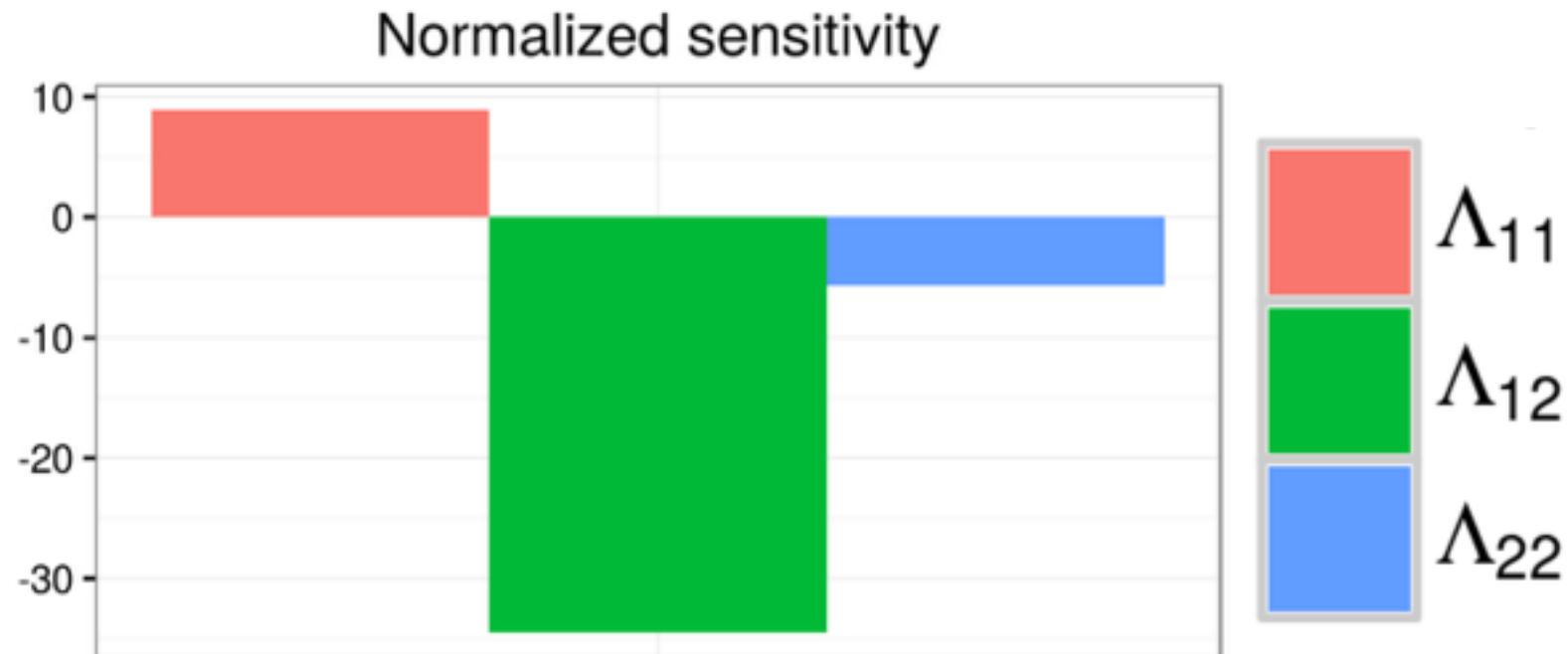
# Microcredit Experiment

- Sensitivity of the expected microcredit effect ( $\tau$ )
- Normalized to be on scale of  $\tau$  std devs
- $\tau$  mean (MFVB): 3.08 USD PPP
- $\tau$  std dev (LRVB): 1.83 USD PPP
- Mean is 1.68 std dev from 0
- $\Lambda_{11} \pm 0.04$



# Microcredit Experiment

- Sensitivity of the expected microcredit effect ( $\tau$ )
- Normalized to be on scale of  $\tau$  std devs
- $\tau$  mean (MFVB): 3.08 USD PPP
- $\tau$  std dev (LRVB): 1.83 USD PPP
- Mean is 1.68 std dev from 0
- $\Lambda_{11} \pm 0.04$   
 $\Rightarrow$  Mean  $> 2$  std dev



# Conclusions

- We provide *linear response variational Bayes*: supplements MFVB for fast & accurate **covariance** estimate
- More from LRVB: fast & accurate **robustness** quantification
- Interested in your data and models:
  - Sensitivity to prior perturbations
  - Sensitivity to likelihood, data perturbations
- Computational statistical trade-offs
  - New data summaries: coresets, approx. sufficient stats
  - Criteo data set: 40 million data points, 3 million features, our runtime: ~20 seconds on 24 cores
  - Theoretical guarantees on finite-sample quality



# References

T Broderick, N Boyd, A Wibisono, AC Wilson, and MI Jordan. Streaming variational Bayes. *NIPS* 2013.

T Campbell\*, JH Huggins\*, J How, and T Broderick. Truncated random measures. Submitted. ArXiv:1603.00861. Poster at ISBA 2016.

**R Giordano, T Broderick, and MI Jordan. Linear response methods for accurate covariance estimates from mean field variational Bayes. *NIPS*, 2015.**

**R Giordano, T Broderick, R Meager, JH Huggins, and MI Jordan. Fast robustness quantification with variational Bayes. *ICML Workshop on #Data4Good: Machine Learning in Social Good Applications*, 2016. ArXiv:1606.07153.**

JH Huggins, T Campbell, and T Broderick. Core sets for scalable Bayesian logistic regression. *NIPS* 2016.

R Meager. Understanding the impact of microcredit expansions: A Bayesian hierarchical analysis of 7 randomised experiments. ArXiv:1506.06669, 2016.

# References

- R Bardenet, A Doucet, and C Holmes. On Markov chain Monte Carlo methods for tall data. arXiv, 2015.
- CM Bishop. *Pattern Recognition and Machine Learning*, 2006.
- D Dunson. Robust and scalable approach to Bayesian inference. Talk at *ISBA* 2014.
- B Fosdick. *Modeling Heterogeneity within and between Matrices and Arrays*, Chapter 4.7. PhD Thesis, University of Washington, 2013.
- DJC MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003.
- M Opper and O Winther. Variational linear response. *NIPS* 2003.
- RE Turner and M Sahani. Two problems with variational expectation maximisation for time-series models. In D Barber, AT Cemgil, and S Chiappa, editors, *Bayesian Time Series Models*, 2011.
- B Wang and M Titterton. Inadequacy of interval estimates corresponding to variational Bayesian approximations. In *AISTATS*, 2004.