



## Overview

Black-box Bayesian inference is hard:

- MCMC can be slow.
- Variational inference can be inaccurate.

**Boosting Variational Inference:**

- **Fast:** optimization-based
- **Nonparametric & Adaptive:** iteratively improves by **adapting to residual**

## Variational Bayes

Variational Bayes approximates true posterior  $p(\theta|X)$  within the closest  $q(\theta)$  within a family of distributions  $\mathcal{H}$ , in terms of discrepancy measure  $\mathcal{D}$  between the two distributions.

$$q^* = \arg \min_{q \in \mathcal{H}} \mathcal{D}(q(\theta), p(\theta|X))$$

**Kullback-Leibler (KL) divergence** is often used as discrepancy measure ( $f = \pi(\theta)p(X|\theta)$ )

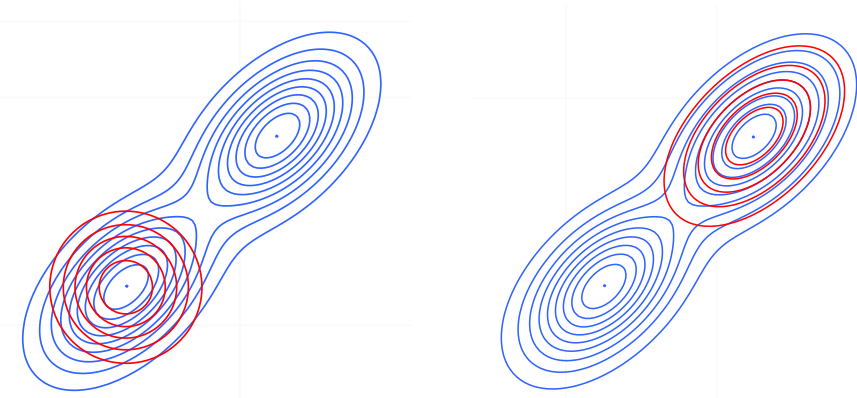
$$\begin{aligned} \mathcal{D}(q, p) &:= \mathcal{D}_{\text{KL}}(q||p) = \int q(\theta) \log \frac{q(\theta)}{p(\theta|X)} d\theta \\ &= \text{const} + \int q \log(q/f) d\theta. \end{aligned}$$

Limitations of current VB:

- Point estimates: **often good, can be biased**
- **Poor uncertainty estimates: covariance, multimodality**
- **Cannot improve accuracy given more time**

## VB Approximation Family

Accuracy of VB is *mainly* limited by the **inflexibility of approximation family**.



1 Mean-field  $q(\theta) = \prod_i q_i(\theta_i)$

2 Full-rank Gaussian

$$\mathcal{H}_1 = \{h : h(\theta) = \mathcal{N}_{\mu, \Sigma}(\theta)\}$$

3 Mixture of  $k$  Gaussians

$$\mathcal{H}_k = \{h : h(\theta) = \sum_{j=1}^k w_j \mathcal{N}_{\mu_j, \Sigma_j}(\theta), \mathbf{w} \in \Delta_k\}$$

4 **Our choice:** All finite Gaussian mixtures

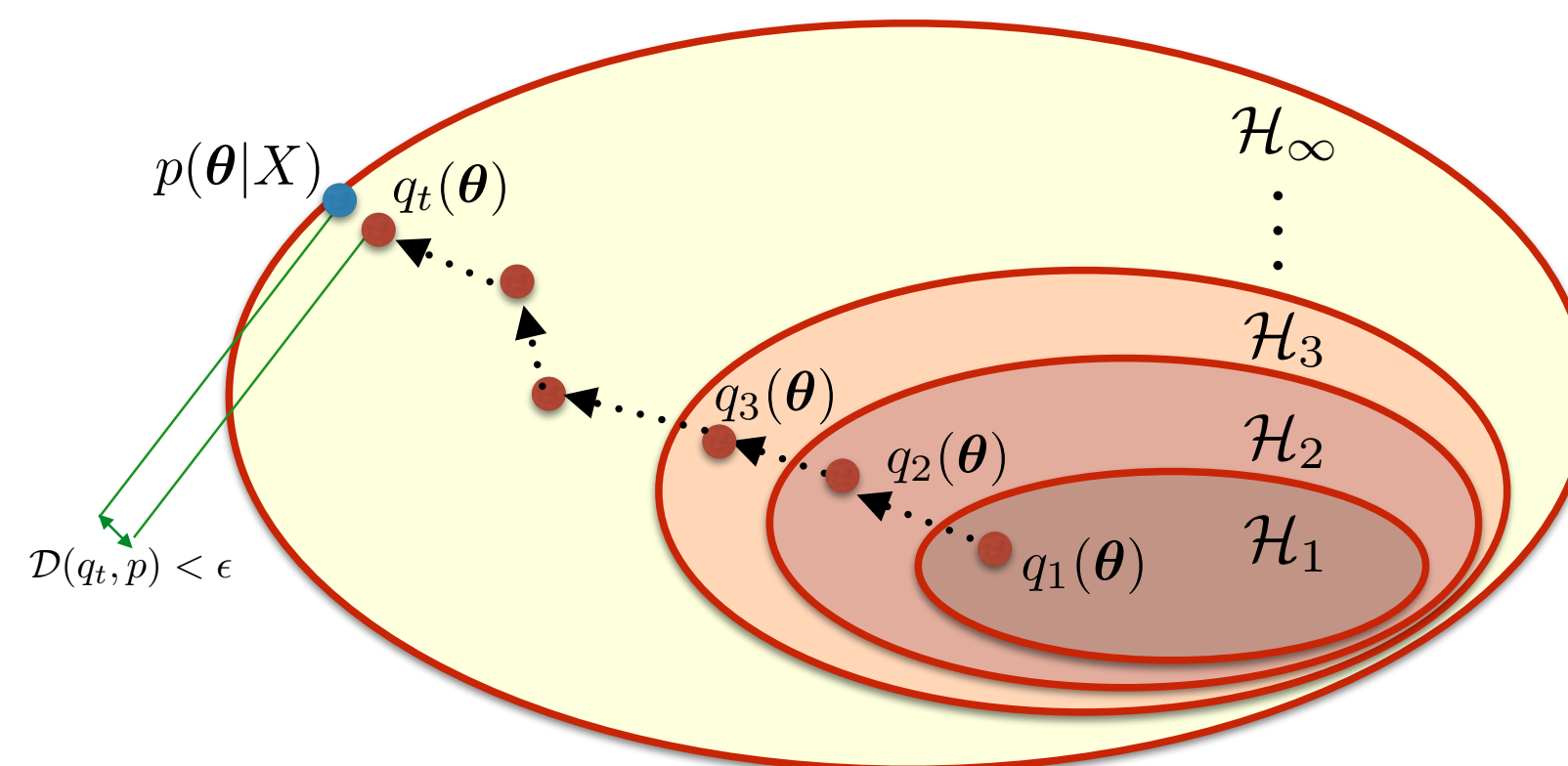
$$\mathcal{H}_{\infty} = \bigcup_{k=1}^{\infty} \mathcal{H}_k$$

Family	Covariance	Multimodality	Arbitrary approximation
Mean-field	✗	✗	✗
Full-rank $\mathcal{N}_{\mu, \Sigma}$	✓	✗	✗
$\mathcal{H}_k$	✓	✓	✗
$\mathcal{H}_{\infty}$	✓	✓	✓

## Greedy Boosting

Want to construct a sequence of approximations  $q_t(\theta) \in \mathcal{H}_t$  such that as  $t \rightarrow \infty$

$$\Delta \mathcal{D}(q_t) := \mathcal{D}(q_t, p) - \inf_{q \in \mathcal{H}_{\infty}} \mathcal{D}(q, p) \searrow 0.$$



## Greedy Boosting Algorithm

1 Start with  $q_1 \in \mathcal{H}_1$ .

2 Then iteratively for  $t = 2, 3, \dots$ , we let

$$q_t = (1 - \alpha_t) q_{t-1} + \alpha_t h_t$$

such that for some  $\epsilon_t \searrow 0$ ,

$$\mathcal{D}(q_t, p) \leq \inf_{h \in \mathcal{H}_t, 0 \leq \alpha \leq 1} \mathcal{D}((1 - \alpha)q_{t-1} + \alpha h, p) + \epsilon_t. (*)$$

**However, optimization (\*) is non-convex.**

## Our Algorithm

**Two-step approach for Greedy Boosting (\*).**

**Step 1: Gradient Boosting: Dist.  $h_t$**

Friedman, (2001) proposed identifying the form of  $h_t$  with the **gradient information** when **increment is small**.

For  $\mathcal{D}_{\text{KL}}$ , the **negative functional gradient** is the **residual** of log posterior density:

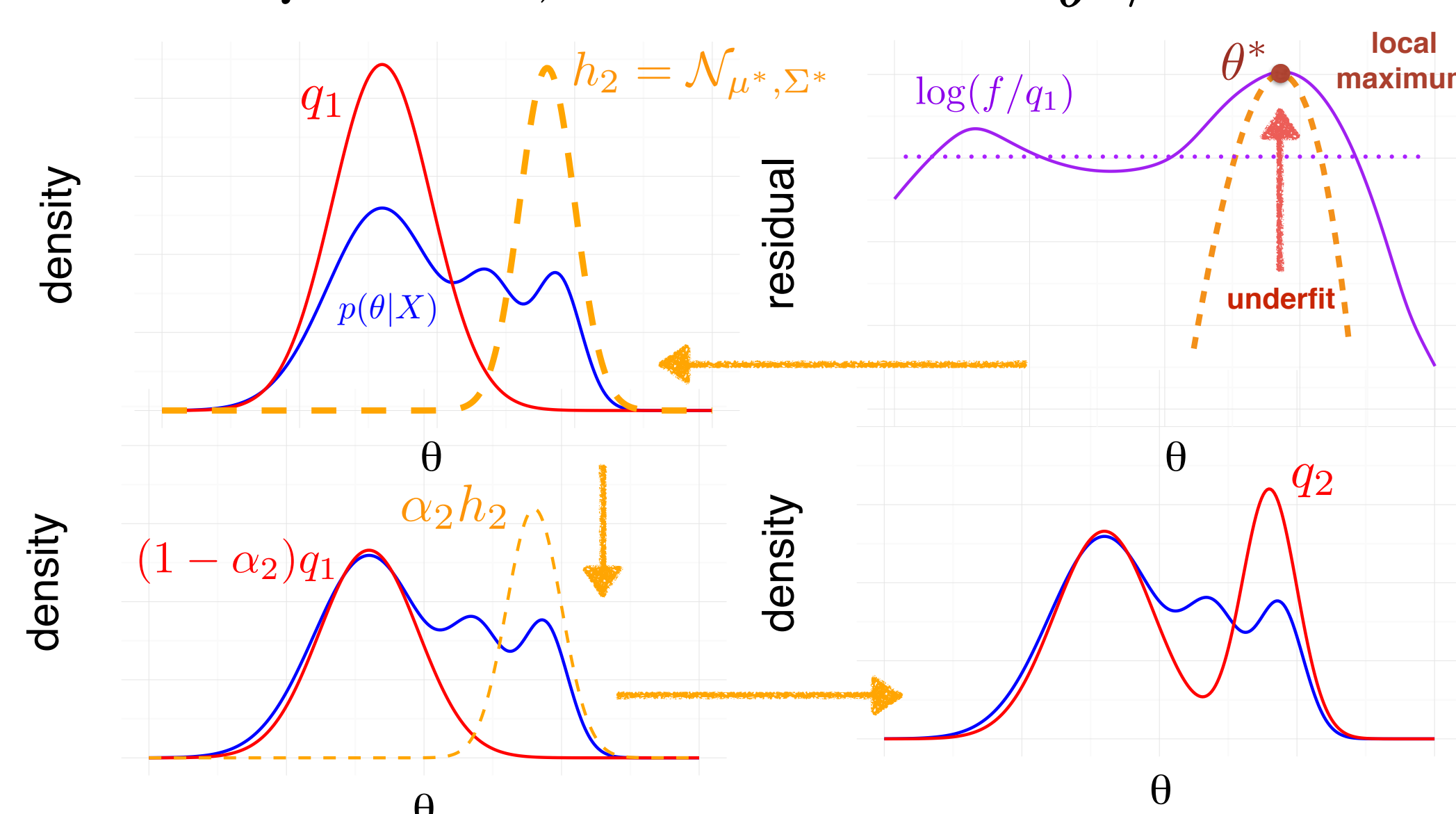
$$-\nabla \tilde{\mathcal{D}}_{\text{KL}}(q_{t-1}) = \log(f(\theta)/q_{t-1}(\theta)).$$

To minimize KL, we match  $h_t$  to  $-\nabla \tilde{\mathcal{D}}_{\text{KL}}(q_{t-1})$ :

$$\hat{h}_t = \arg \min_{h \in \mathcal{H}_t, c > 0} \|c \cdot h - \log(f/q_{t-1})\|_2^2.$$

With Laplacian approximation to the residual, we have a simple algorithm for quickly identifying  $\hat{h}_{\mu_t, \Sigma_t}$  using optimization. We have closed-form solutions:

$$\mu^* = \theta^*, \quad \Sigma^* = \text{Hessian}_{\theta^*}^{-1}/2.$$



**Step 2: Stochastic Newton's: Weight  $\alpha_t$**

Fixing  $h_t$ , determining corresponding weight

$$\alpha_t = \min_{0 \leq \alpha \leq 1} \tilde{\mathcal{D}}_{\text{KL}}((1 - \alpha)q_{t-1} + \alpha h_t)$$

is **convex**. Further, by drawing samples from  $q_{t-1}$  and  $h_t$ , we can get **Monte Carlo estimates** of derivatives  $\tilde{\mathcal{D}}'_{\text{KL}}$  and  $\tilde{\mathcal{D}}''_{\text{KL}}$ .

## Theoretical Results

From Zhang, (2003), for greedy boosting, if  $\mathcal{D}(q, p)$  is (1) **convex** in  $q$  and (2) **strongly smooth** in  $q$ , then we have

$$\Delta \mathcal{D}(q_t) \rightarrow 0 \text{ at rate } O(1/t).$$

In **Theorem 1**, we showed that under mild conditions (e.g., that hold on a bounded set)  $\mathcal{D}_{\text{KL}}$  satisfies these conditions.

## Simulation Experiments

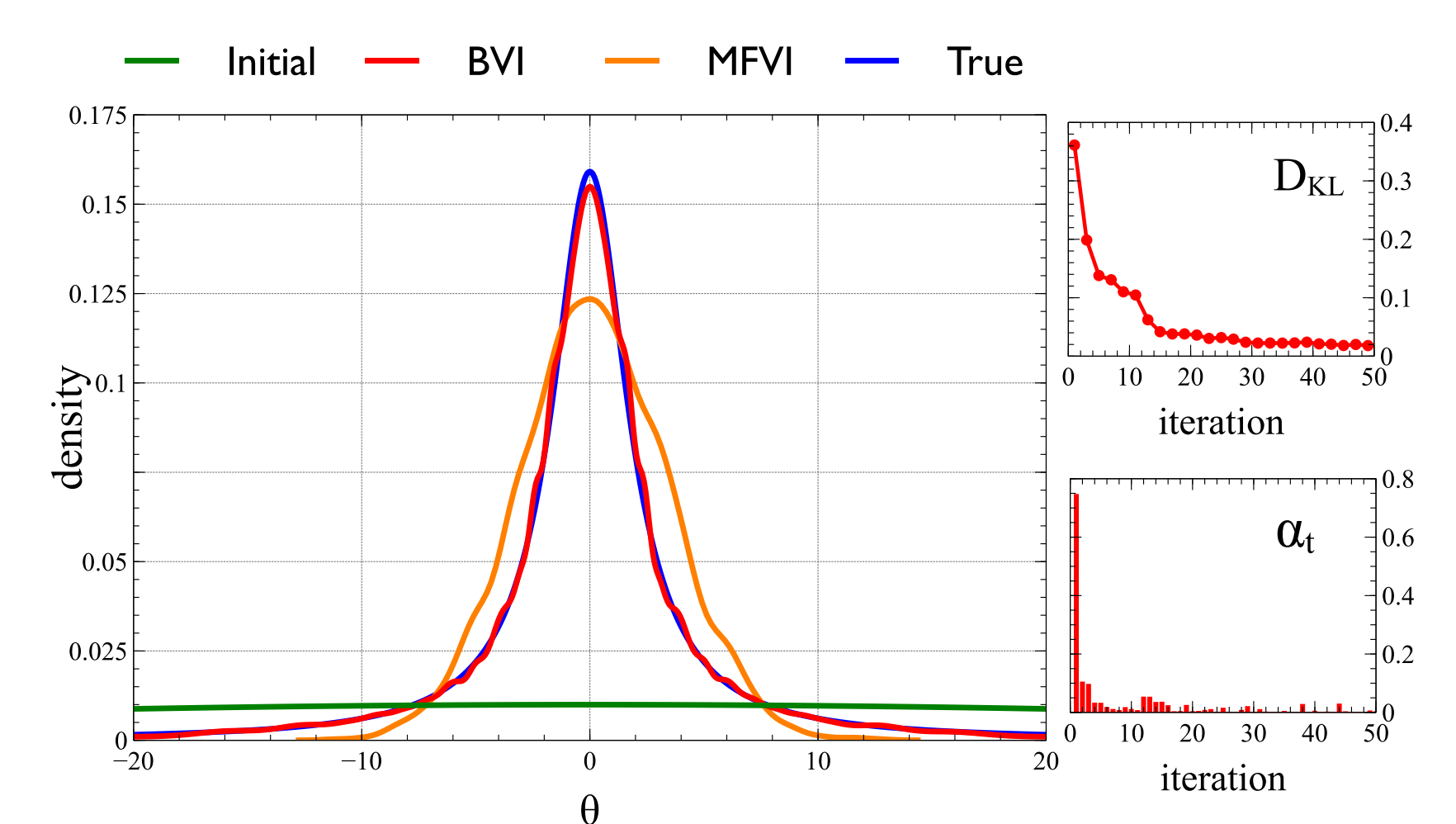


Figure 1: True: Heavy-tailed Cauchy

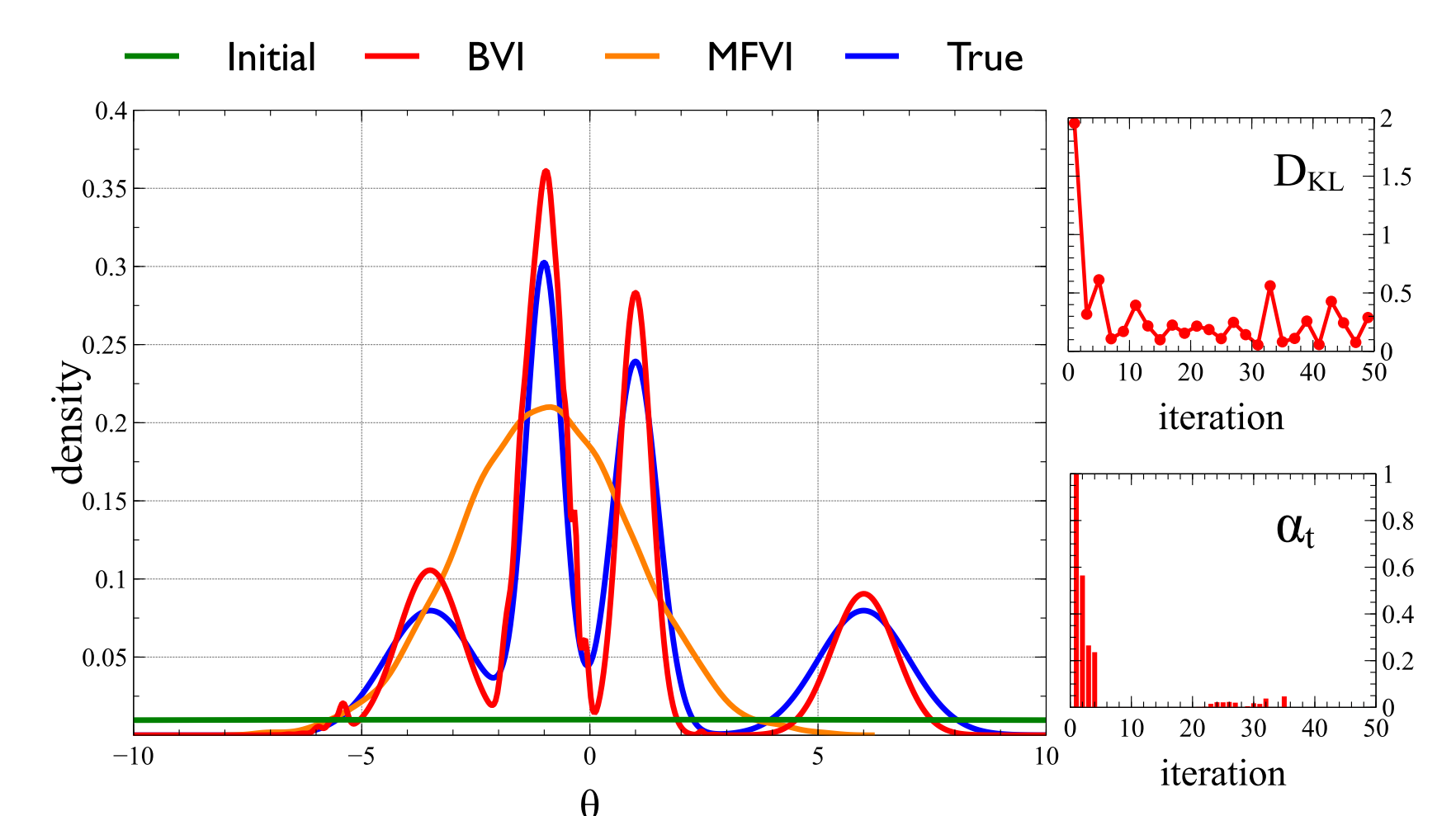


Figure 2: True: Mixture of univariate Gaussians

## Logistic Regression Experiment

We run Bayesian logistic regression on the Noda1 dataset, consisting of  $N = 53$  observations of  $d = 6$  predictors  $x_i$  and a binary response  $y_i \in \{-1, +1\}$ . We compare to MCMC (as truth) and mean-field VB.

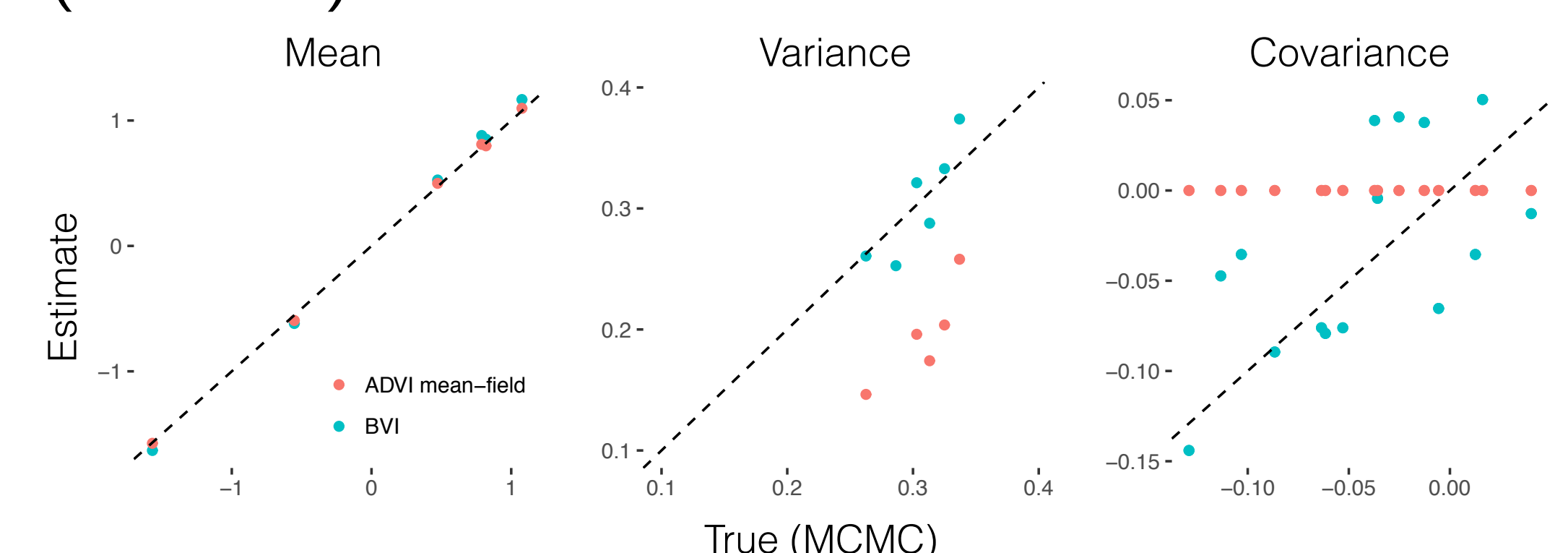


Figure 3: Bayesian logistic regression

## References

- Friedman, Jerome H (2001). "Greedy function approximation: a gradient boosting machine". In: *Annals of statistics*, pp. 1189–1232.
- Zhang, Tong (2003). "Sequential greedy approximation for certain convex optimization problems". In: *IEEE Transactions on Information Theory* 49.3, pp. 682–691.