



# Generalized Variational Inference (GVI)

Posterior beliefs with the rule of three

---

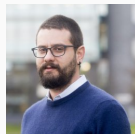
Jeremias Knoblauch<sup>1,3</sup>, Jack Jewson<sup>1,3</sup>,  
Theodoros Damoulas<sup>1,2,3</sup>

September 29, 2020

<sup>1</sup>University of Warwick, Department of Statistics

<sup>2</sup>University of Warwick, Department of Computer Science

<sup>3</sup>The Alan Turing Institute for Data Science and AI



# Structure of the talk

1. The form of the Bayesian problem
  - 1.1 The traditional perspective
  - 1.2 The optimization perspective
  - 1.3 The loss-minimization perspective
  - 1.4 The new perspective
2. The form of the Generalized Bayesian problem
  - 2.1 Provable modularity
  - 2.2 Axiomatic derivation
  - 2.3 Relationship to existing methods
3. Reinterpreting standard **VI**
  - 3.1 Optimality & reinterpretation of standard **vi**
  - 3.2 Why does **dvi** produce better posteriors?
  - 3.3 Towards **gvi**
4. **GVI**: What does it do?
  - 4.1 M-open vs M-closed
  - 4.2 The losses
  - 4.3 Uncertainty Quantification
  - 4.4 Three **gvi** use cases
5. **GVI**: Inference & Experiments
  - 5.1 Black box **gvi**
  - 5.2 Robust Bayesian On-line Changepoint Detection
  - 5.3 Bayesian Neural Networks
  - 5.4 Deep Gaussian Processes

# 1 The form of the Bayesian problem

**Purpose of part 1:** Motivate the rule of three

- (1) Bayesian inference minimizes **losses**
- (2) Bayesian inference **regularizes** with the prior
- (3) Bayesian inference = **optimization** over **(sub)spaces of probability measures**

# 1.1 The Bayesian problem: Traditional perspective

**Ingredients** (for the simplest case) are:

- $n = n_1 + n_2$  observations  $\mathbf{x} = (x_1, x_2, \dots, x_{n_1+n_2})^T$ ,
- **prior**  $\pi(\boldsymbol{\theta})$ ,
- **likelihoods**  $\{p(x_i|\boldsymbol{\theta})\}_{i=1}^{n_1+n_2}$

**Output = posterior belief:**

$$q^*(\boldsymbol{\theta}) \propto \pi(\boldsymbol{\theta}) \prod_{i=1}^{n_1+n_2} p(x_i|\boldsymbol{\theta}) = \tilde{\pi}(\boldsymbol{\theta}) \prod_{i=n_1+1}^{n_2} p(x_i|\boldsymbol{\theta}), \text{ for } \tilde{\pi}(\boldsymbol{\theta}) = \pi(\boldsymbol{\theta}) \prod_{i=1}^{n_1} p(x_i|\boldsymbol{\theta})$$

**Inference interpretation = belief updates:**

- likelihoods  $\{p(x_i|\boldsymbol{\theta})\}_{i=1}^{n_1+n_2}$  update prior about  $\boldsymbol{\theta}$
- Old posterior  $\tilde{\pi}(\boldsymbol{\theta}) =$  new prior (**coherence/Bayesian additivity**)

## 1.2 The Bayesian problem: The optimization perspective

Zellner (1988) shows that the Bayes posterior  $q^*(\theta)$  solves

$$q^*(\theta) = \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \underbrace{\mathbb{E}_{q(\theta)} \left[ \sum_{i=1}^n -\log(p(x_i|\theta)) \right]}_{\text{minimized by } q(\theta) = \delta_{\hat{\theta}_n}(\theta), \hat{\theta}_n = \text{MLE}} + \underbrace{\text{KLD}(q||\pi)}_{\text{minimized by } q = \pi} \right\}, \quad (1)$$

### Notation:

- $\mathcal{P}(\Theta)$  = all probability distributions on  $\Theta$
- **KLD** = Kullback-Leibler divergence =  $\mathbb{E}_{q(\theta)} [\log q(\theta) - \log \pi(\theta)]$

### Inference interpretation = regularized loss-minimization:

- $-\log(p(x_i|\theta))$  = **loss** of  $\theta$  for  $x_i$
- Inference = regularizing MLE  $\hat{\theta}_n$  with **KLD**( $q||\pi$ )

## 1.3 The Bayesian problem: The loss-minimization perspective

Bissiri et al. (2016): Bayes posteriors  $q^*(\theta)$  for general loss  $\ell(\theta, x_i)$ :

$$q^*(\theta) \propto \pi(\theta) \exp \left\{ - \sum_{i=1}^{n_1+n_2} \ell(\theta, x_i) \right\} = \tilde{\pi}(\theta) \exp \left\{ - \sum_{i=n_1+1}^{n_2} \ell(\theta, x_i) \right\}$$

for  $\tilde{\pi}(\theta) = \pi(\theta) \exp \left\{ - \sum_{i=1}^{n_1} \ell(\theta, x_i) \right\}$

**Inference interpretation = belief updates:**

- Again: losses  $\{\ell(\theta, x_i)\}_{i=1}^{n_1+n_2}$  update prior about  $\theta$
- Again: Old posterior  $\tilde{\pi}(\theta)$  = new prior (**coherence**)
- Difference:  $\theta$  arbitrary, e.g.  $\ell(\theta, x_i) = |x_i - \theta|$  admissible

## 1.4 The Bayesian problem: The new perspective I/II

**Easy to show:** Zellner's representation valid for any  $\ell(\boldsymbol{\theta}, x_i)$ :

$$q^*(\boldsymbol{\theta}) = \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \underbrace{\mathbb{E}_{q(\boldsymbol{\theta})} \left[ \sum_{i=1}^n \ell(\boldsymbol{\theta}, x_i) \right]}_{\text{minimized by } \delta_{\hat{\boldsymbol{\theta}}_n}(\boldsymbol{\theta})} + \underbrace{\text{KLD}(q \parallel \pi)}_{\text{minimized by } q = \pi} \right\}$$

**Bissiri et al. (2016)'s generalization** (preserves coherence):

- Replacing  $-\log(p(x_i|\boldsymbol{\theta}))$  with other losses  $\ell(\boldsymbol{\theta}, x_i)$

**Two more generalizations** (break coherence):

- Replacing  $\mathcal{P}(\Theta)$  with  $\mathcal{Q} \subset \mathcal{P}(\Theta)$  (= VI)
- Replacing **KLD** with inference-problem specific regularizers

## 1.4 The Bayesian problem: The new perspective II/II

Our generalized representation of Bayesian inference:

$$q^*(\theta) = \arg \min_{q \in \Pi} \left\{ \underbrace{\mathbb{E}_{q(\theta)} \left[ \sum_{i=1}^n \ell(\theta, x_i) \right]}_{\text{minimized by } \delta_{\hat{\theta}_n}(\theta)} + \underbrace{D(q||\pi)}_{\text{minimized by } q = \pi} \right\}$$

**Notation:**

- if  $\Pi$  = variational family, write  $\mathcal{Q}$ .
- $\ell_n(\theta, \mathbf{x}) = \sum_{i=1}^n \ell(\theta, x_i)$

**Inference interpretation = regularized & constrained minimization:**

- $\ell_n(\theta, \mathbf{x})$  = **loss** of  $\theta$  to minimize
- **D** = **divergence**, acting as uncertainty quantifier/regularizer
- $\Pi$  = set of **admissible posterior** beliefs
- **Inference** = constrained, regularized optimization

⇒ **Shorthand Notation:**  $P(\ell_n, D, \Pi)$



## 2 The form of the Generalized Bayesian problem

**Purpose of part 2:** Investigate  $P(\ell_n, D, \Pi)$

- (1) Interpretations & modularity of  $\ell_n$ ,  $D$  and  $\Pi$ ?
- (2) Is there an axiomatic justification?
- (3) Which existing methods does this (not) encompass?

## 2.1 Generalized Bayesian problem: provable modularity

$$q^*(\theta) = \arg \min_{q \in \Pi} \left\{ \underbrace{\mathbb{E}_{q(\theta)} \left[ \sum_{i=1}^n \ell(\theta, x_i) \right]}_{\text{minimized by } \delta_{\hat{\theta}_n}(\theta)} + \underbrace{D(q \parallel \pi)}_{\text{minimized by } q = \pi} \right\}$$

Roles of  $\ell_n$ ,  $D$ ,  $\Pi$ :

- $\ell_n$ : which parameter  $\theta$  do we care about?
  - $D$ : How is uncertainty quantified/what does  $q^*$  look like?
  - $\Pi$ : Which beliefs are allowed?
- $\Rightarrow$  (provable) modularity of  $P(\ell_n, D, \Pi)$ !

### Theorem 1 (GVI modularity)

For Bayesian inference with  $P(\ell_n, D, \Pi)$ , making it robust to model misspecification amounts to changing  $\ell_n$ . Conversely, adapting uncertainty quantification (fixing  $\Pi$ ,  $\pi$ ,  $\theta^*$ ,  $\hat{\theta}_n$ ) amounts to changing  $D$ .

## 2.2 Generalized Bayesian problem: Axiomatic derivation I/II

### Axiom 1 (Representation)

The posterior  $q^* \in \mathcal{P}(\Theta)$  is constructed by solving an optimization problem over some space  $\Pi \subseteq \mathcal{P}(\Theta)$ . The optimization seeks to minimize exactly two criteria that do not interact:

- (i) The in-sample loss  $\sum_{i=1}^n \ell(\theta, x_i)$  to be expected under  $q^*(\theta)$ .
- (ii) The deviation from the prior  $\pi(\theta)$  as measured by the magnitude of some statistical divergence  $D$ .

### Axiom 2 (Information Equivalence)

Take any two constants  $C, M \in \mathbb{R}_{>0}$  and let the posteriors  $q_1^*, q_2^* \in \mathcal{P}(\Theta)$  be computed based on the same optimization problem, but with two different loss functions  $\ell^{(1)}$  and  $\ell^{(2)}$  as well as two different divergences  $D^{(1)}, D^{(2)}$  so that  $D^{(1)} = D^{(2)} + M$  and so that  $\ell^{(1)} = \ell^{(2)} + C$ . Then,  $q_1^*(\theta) = q_2^*(\theta)$  (almost everywhere).

## 2.2 Generalized Bayesian problem: Axiomatic derivation I/II

### Axiom 3 (Generalized Likelihood Principle)

Suppose the posteriors  $q_n^*, q_{n+m}^* \in \Pi$  are computed based on an optimization problem satisfying Axiom 1. Assume that the same optimization problem is used for both posteriors, except for a difference in the samples  $x_{1:n}$  and  $x_{1:n+m}$  and potentially different loss functions  $\ell^{(1)}$  and  $\ell^{(2)}$  that are used, respectively.

- (i) Provided that there is an information difference between  $x_{1:n}$  and  $x_{1:n+m}$  for  $m > 0$ , the posteriors are different. In other words,  $\sum_{i=1}^n \ell^{(1)}(\theta, x_i) \neq \sum_{i=1}^{n+m} \ell^{(2)}(\theta, x_i)$  implies that  $q_n^* \neq q_{n+m}^*$ , even if  $\ell^{(1)} = \ell^{(2)}$ .
- (ii) Provided that there is a difference in the measure of information between  $\ell^{(1)}$  and  $\ell^{(2)}$ , the posteriors are different. In other words,  $\ell^{(1)} \neq \ell^{(2)}$  implies that  $q_n^* \neq q_{n+m}^*$ , even if  $m = 0$  so that  $x_{1:n} = x_{1:n+m}$ .

## 2.2 Generalized Bayesian problem: Axiomatic derivation II/II

### Theorem 2 (Form 1)

If Axiom 1 holds,  $P(\ell_n, D, \Pi)$  has form  $\arg \min_{q \in \Pi} \{L(q|\mathbf{x}, \ell_n, D)\}$  for  $L(q|\mathbf{x}, \ell_n, D) = f(\mathbb{E}_{q(\theta)}[\ell_n(\theta, \mathbf{x})], D(q|\pi))$ , for some  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ .

### Theorem 3 (Form 2)

For  $P(\ell_n, D, \Pi)$  being  $\arg \min_{q \in \Pi} \{L(q|\mathbf{x}, \ell_n, D)\}$  and  $\circ$  an elementary operation on  $\mathbb{R}$ ,  $L(q|\mathbf{x}, \ell_n, D) = \mathbb{E}_{q(\theta)}[\ell_n(\theta, \mathbf{x})] \circ D(q|\pi)$  satisfies Axioms 1, 2 and 3 only if  $\circ = +$ .

### Implications/relevance:

- Bayesian inference = constrained, regularized optimization
- Objective only depends on  $\mathbb{E}_{q(\theta)}[\ell_n(\theta, \mathbf{x})]$  and  $D(q|\pi)$
- For elementary  $f(\mathbb{E}_{q(\theta)}[\ell_n(\theta, \mathbf{x})], D(q|\pi))$ ,  $f$  must be addition.  
(**Note:** Axiom 2 excludes most non-elementary  $f$ )

## 2.3 Generalized Bayesian problem & existing methods I/III

$$q^*(\theta) = \arg \min_{q \in \Pi} \left\{ \mathbb{E}_{q(\theta)} [\ell_n(\theta, \mathbf{x})] + D(q \parallel \pi) \right\}$$

$P(\ell_n, D, \Pi)$  covers & gives **insight** into existing methods, e.g.

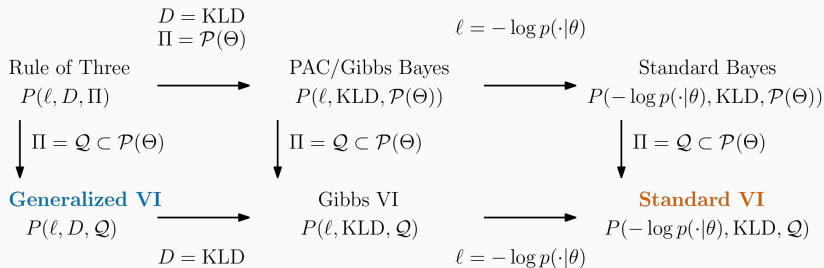
- **Power Bayes:**  $P(w\ell_n, D, \Pi) = P(\ell_n, \frac{1}{w}D, \Pi)$ .  
( $\implies$   $w$ -power likelihood =  $\frac{1}{w} \times$  **more trust in your prior.**)
- **Regularized Bayes:** Adding  $\Phi(q(\theta, \mathbf{x})) = \mathbb{E}_{q(\theta, \mathbf{x})} [\phi(\theta, \mathbf{x})]$  into the objective corresponds to  $P(\ell_n + \phi, D, \Pi)$ .  
( $\implies$  RegBayes = a form of **GVI** that changes  $\ell_n$ )
- **PAC Bayes:** Operationalizes a variational inference procedure based on PAC-bounds whose complexity penalty is some  $D \neq \text{KLD}$  (e.g. Alquier and Guedj, 2018)

## 2.3 Generalized Bayesian problem & existing methods II/III

Method	$\ell(\boldsymbol{\theta}, \mathbf{x}_i)$	$D$	$\Pi$
Standard Bayes	$-\log(p(\boldsymbol{\theta} \mathbf{x}_i))$	KLD	$\mathcal{P}(\Theta)$
Generalized Bayes <sup>1</sup>	any $\ell$	KLD	$\mathcal{P}(\Theta)$
Power Bayes <sup>2</sup>	$-\log(p(\boldsymbol{\theta} \mathbf{x}_i))$	$\frac{1}{w}$ KLD, $w > 1$	$\mathcal{P}(\Theta)$
Divergence Bayes <sup>3</sup>	divergence-based $\ell$	KLD	$\mathcal{P}(\Theta)$
<b>Standard VI</b>	<b><math>-\log(p(\boldsymbol{\theta} \mathbf{x}_i))</math></b>	<b>KLD</b>	$\mathcal{Q}$
Power VI <sup>4</sup>	$-\log(p(\boldsymbol{\theta} \mathbf{x}_i))$	$\frac{1}{w}$ KLD, $w > 1$	$\mathcal{Q}$
Regularized Bayes <sup>5</sup>	$-\log(p(\boldsymbol{\theta} \mathbf{x}_i)) + \phi(\boldsymbol{\theta}, \mathbf{x}_i)$	KLD	$\mathcal{Q}$
Gibbs VI <sup>6</sup>	any $\ell$	KLD	$\mathcal{Q}$
<b>Generalized VI</b>	<b>any <math>\ell</math></b>	<b>any <math>D</math></b>	$\mathcal{Q}$

**Table 1** –  $P(\ell_n, D, Q)$  & existing methods. <sup>1</sup>(Bissiri et al., 2016), <sup>2</sup>(e.g. Holmes and Walker, 2017; Grünwald et al., 2017; Miller and Dunson, 2018), <sup>3</sup>(e.g. Hooker and Vidyashankar, 2014; Ghosh and Basu, 2016; Futami et al., 2017; Jewson et al., 2018), <sup>4</sup>(e.g. Yang et al., 2017; Huang et al., 2018) <sup>5</sup>(Ganchev et al., 2010; Zhu et al., 2014), <sup>6</sup>(Alquier et al., 2016; Futami et al., 2017)

## 2.3 Generalized Bayesian problem & existing methods II/III



**Figure 1** – A taxonomy of some important belief distributions as special cases.



## 2.3 Generalized Bayesian problem & existing methods III/III

Not everything fits  $P(\ell_n, D, \Pi)$ :

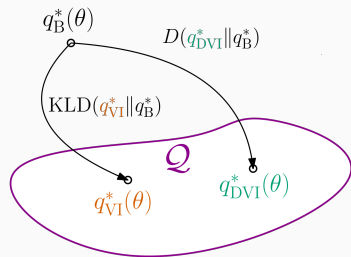
- (1) **Laplace approximations** (e.g., INLA)
- (2) **Divergence-Variational inference (DVI)**: VI based on discrepancy  $D \neq \text{KLD}$  (locally) solving  $q^* = \arg \min_{q \in \mathcal{Q}} D(q \parallel \tilde{q})$  for  $\tilde{q} =$  standard Bayesian posterior, e.g.
  - $D = \text{Rényi's } \alpha\text{-divergence}$  (Li and Turner, 2016; Saha et al., 2017)
  - $D = \chi\text{-divergence}$  (Dieng et al., 2017)
  - $D = \text{operators}$  (Ranganath et al., 2016)
  - $D = \text{scaled AB-divergence}$  (Regli and Silva, 2018)
  - $D = \text{Wasserstein distance}$  (Ambrogioni et al., 2018)
  - ...
- (3) **Expectation Propagation (EP)** (Minka, 2001; Opper and Winther, 2000) and its variants (e.g. Hernández-Lobato et al., 2016).  
**Note:** Particular type of **DVI**, with  $D =$  (local) reverse KLD

### Purpose of part 3: Motivating GVI

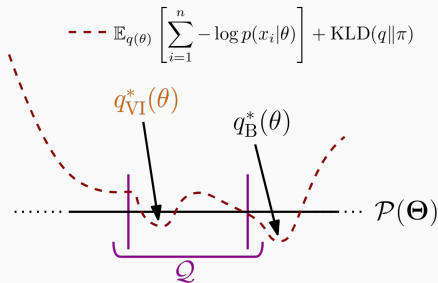
- (1) **Standard VI**: Optimality & reinterpretation
- (2) **DVI**: “suboptimal” methods with better posteriors
- (3) **GVI**: A modular alternative to **DVI**

## 3.1 Optimality & reinterpretation of standard VI I/V

Relationship between VI, GVI, DVI and exact inference?



(a) DVI (projection-centric)  
interpretation of VI



(b) GVI (= optimization-centric)  
Interpretation of VI

## 3.1 Optimality & reinterpretation of standard VI II/V

**Standard VI:**  $q_{VI}^* = \arg \min_{q \in \mathcal{Q}} \text{KLD}(q \| q_B^*)$ ,  $q_B^*$  solves

$$P(\ell_n, \text{KLD}, \mathcal{P}(\Theta))$$

$$\text{KLD}(q \| q_B^*) = \underbrace{\mathbb{E}_{q(\theta)} \left[ \log \left( \frac{q(\theta)}{\exp \left\{ - \sum_{i=1}^n \ell(\theta, x_i) \right\} \pi(\theta)} \right) \right]}_{\text{(Generalized) ELBO}} + \underbrace{\log \left( \int_{\theta} \exp \left\{ - \sum_{i=1}^n \ell(\theta, x_i) \right\} \pi(\theta) d\theta \right)}_{\text{Generalized 'log evidence'}}$$

Inference = minimizing ELBO, which you can rewrite as

$$\text{ELBO}(q) = \mathbb{E}_{q(\theta)} \left[ \sum_{i=1}^n \ell(\theta, x_i) \right] + \text{KLD}(q \| \pi). \quad (2)$$

... which is exactly the objective of  $P(\ell_n, \text{KLD}, \mathcal{Q})$

## 3.1 Optimality & reinterpretation of standard VI III/V

In other words,  $P(\ell_n, \text{KLD}, \mathcal{Q})$  (= ELBO) is

$$q_{\text{VI}}^*(\theta) = \arg \min_{q \in \mathcal{Q}} \left\{ \mathbb{E}_{q(\theta)} [\ell_n(\theta, \mathbf{x})] + \text{KLD}(q || \pi) \right\},$$

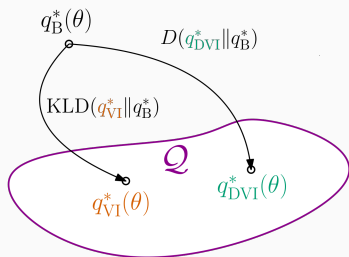
the  $\mathcal{Q}$ -constrained relaxation of  $P(\ell_n, \text{KLD}, \mathcal{P}(\Theta))$ , whose objective is

$$q_B^*(\theta) = \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \mathbb{E}_{q(\theta)} [\ell_n(\theta, \mathbf{x})] + \text{KLD}(q || \pi) \right\},$$

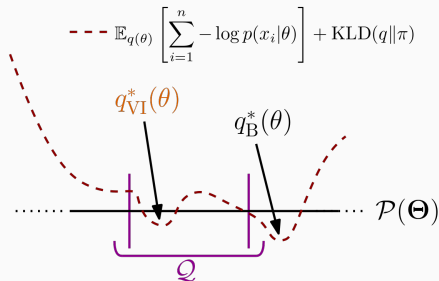
(which is the exact Bayesian objective).

⇒ Reinterpretation of **standard VI** as **Constrained optimization!**

## 3.1 Optimality & reinterpretation of standard VI IV/V



(a) **DVI** (projection-centric)  
interpretation of **VI**



(b) **GVI** (= optimization-centric)  
Interpretation of **VI**

## 3.1 Optimality & reinterpretation of standard VI V/V

**Consequence I/II: VI-optimality**

**Theorem 4 (VI optimality)**

For exact and coherent Bayesian posteriors solving  $P(\ell_n, \text{KLD}, \mathcal{P}(\Theta))$  and a fixed variational family  $\mathcal{Q}$ , standard VI produces the uniquely optimal  $\mathcal{Q}$ -constrained approximation to  $P(\ell_n, \text{KLD}, \mathcal{P}(\Theta))$ . Having decided on approximating the Bayesian posterior with some  $q \in \mathcal{Q}$ , VI provides the uniquely optimal solution.

## 3.2 Why does DVI produce better posteriors? I/II

**Consequences II/II:** DVI-suboptimality. **Three** big disadvantages:

- (1) If  $F \neq \text{KLD}$ , DVI violates Axioms 1–3.
- (2) DVI conflates  $\ell_n$  and  $D$  (i.e., modularity of  $P(\ell_n, D, \Pi)$  lost).
- (3) Last Thm: DVI gives worse  $\mathcal{Q}$ -constrained posterior than standard VI  
(relative to the standard Bayesian problem  $P(\ell_n, \text{KLD}, \mathcal{P}(\Theta))$ )

**Objection!** DVI often produces better posteriors than **standard VI!**



## 3.2 Why does DVI produce better posteriors? II/II

### Seeming contradiction:

- (1) VI is the best approximation to the *standard* Bayesian posterior
- (2) DVI often outperforms VI (e.g., on test scores)

### Resolution:

DVI outperforms VI by implicitly **defining a new, non-standard Bayesian problem\***

⇒ Inspires **Generalized Variational Inference (GVI)**

\***non-standard Bayesian problem** = problem defining a belief distribution that can be understood as updating a prior belief with some data, but which is **not** obtained by Bayes' rule

## 3.3 Towards GVI I/II

**GVI** = combining advantages of **VI** and **DVI**:

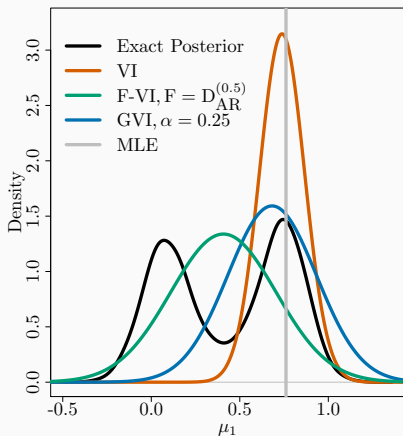
- (1) Has form  $P(\ell_n, D, Q)$  Like **VI** i.e.
  - (i) satisfies Axioms 1–3;
  - (ii) provably interpretable modularity (**loss, uncertainty quantifier, admissible posteriors**)
- (2) Derives different & more appropriate posteriors like **DVI** but
  - (i) without conflating  $\ell_n$  and  $D$
  - (ii) with **explicit** rather than **implicit** changes .

### Definition 1 (GVI)

Any Bayesian inference method solving  $P(\ell_n, D, Q)$  with admissible choices  $\ell_n$ ,  $D$  and  $Q$  is a Generalized Variational Inference (**GVI**) method (and satisfies Axioms 1–3 by definition).

### 3.3 Towards GVI II/II

Illustration: **DVI** aims for  $D$ , but changes  $\ell_n$  – **GVI** doesn't



**Figure 4** – Exact, **VI**, **DVI** ( $F = D_{AR}^{(0.5)}$ ) and  $P(\ell_n, D_{AR}^{(\alpha)}, \mathcal{Q})$  based **GVI** marginals of the location in a 2 component mixture model. Respecting  $\ell_n$ , **VI** and **GVI** provide uncertainty quantification around the most likely value  $\hat{\theta}_n$  via  $D$ . In contrast, **DVI** implicitly changes the loss and has a mode at the locally most *unlikely* value of  $\theta$ .

## 4 GVI: What does it do?

**Purpose of part 4:** Exploring **GVI**'s relationship to M-open world & study three use cases

- (1) Embed **GVI** into the M-open world
- (2) Robust alternatives to  $\ell(\boldsymbol{\theta}, x_i) = -\log(p(x_i|\boldsymbol{\theta}))$
- (3) Prior-robust uncertainty quantification and adjusting marginal variances via  $D$

## 4.1 M-open vs M-closed inference

**M-closed** view: Model and prior are correct, i.e.

$$\exists \theta^* \in \Theta \text{ s.t. } x_i \sim p(x_i | \theta^*)$$

**Q:** Why do purist Bayesians love **exact** inference (MCMC, SMC, ...)?

**A:** With M-closed view, can focus on computation of Bayes posterior

$$q^*(\theta) \propto \prod_{i=1}^n p(x_i | \theta) \pi(\theta) = \text{solution of } P \left( \sum_{i=1}^n -\log p(x_i | \theta), D, \mathcal{P}(\Theta) \right)$$

**M-open** view: Model and prior are **not** correct, i.e.

$$\nexists \theta^* \in \Theta \text{ s.t. } x_i \sim p(x_i | \theta^*)$$

**Traditional stats:** Devise better models until they are 'close enough'

⇒ Focus on inference still useful!

**Modern stats/ML:** Use some black box model (BNN, DGP, ...)

⇒ Do we even want standard posterior!?

## 4.1 Consequences for inference

**Conclusion 1:** If your model is pretty good, use exact inference or **VI**.

**Conclusion 2:** If **DVI** works better than **VI**, your model is not great.

**Conclusion 3:** If you know why model isn't great, address it with **GVI**.

**M-closed** assumptions are very far from the truth sometimes:

- Time Series/on-line inference:
  - (i) Stationarity
  - (ii) Short-term memory
  - (iii) Noise level constant/non-erratic
- Bayesian Neural Networks (BNNs):
  - (i) There is a true weight parameter  $\theta^*$  indexing the network
  - (ii) All priors on the thousands of entries of  $\theta$  is well-specified
- (Deep) Gaussian Processes (DGPs):
  - (i) (Conditionally on latent space) correct likelihood is specified
  - (ii) DGP prior and kernel choice generate correct latent spaces

⇒ Clearly, **M-open** world more appropriate ⇒ **GVI**

## 4.2 GVI: The losses I/III

**GVI modularity:** The loss  $\ell_n$

**Q1:** Why use  $\ell_n(\boldsymbol{\theta}, \mathbf{x}) = \sum_{i=1}^n -\log(p(x_i|\boldsymbol{\theta}))$ ?

**A:** Assuming that the true data-generating mechanism is  $\mathbf{x} \sim g$ ,

$$\begin{aligned} \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n -\log(p(x_i|\boldsymbol{\theta})) &\approx \arg \min_{\boldsymbol{\theta}} \mathbb{E}_g [-\log(p(\mathbf{x}|\boldsymbol{\theta}))] \\ &= \arg \min_{\boldsymbol{\theta}} \mathbb{E}_g [-\log(p(\mathbf{x}|\boldsymbol{\theta})) + \log(g(\mathbf{x}))] = \arg \min_{\boldsymbol{\theta}} \text{KLD}(g\|p(\cdot|\boldsymbol{\theta})) \end{aligned}$$

**Interpretation:**  $-\log(p(x_i|\boldsymbol{\theta}))$  = targeting KLD-minimizing  $p(\cdot|\boldsymbol{\theta})$

**Q2:** Are there other  $\mathcal{L}^D(p(x_i|\boldsymbol{\theta}))$  for divergence  $D$ ?

**A:** Yes! (e.g. Jewson et al., 2018; Futami et al., 2017; Ghosh and Basu, 2016; Hooker and Vidyashankar, 2014)

## 4.2 GVI: The losses II/III

**Q3:** Why use other  $\mathcal{L}^D(p(x_i|\theta))$ ?

**A:** Robustness (for  $D$  = a robust divergence) [log/KLD non-robust!]

**Robustness recipe:**  $\alpha/\beta/\gamma$ -divergences using generalized log functions

**E.g.:**  $\beta$  indexes  $\beta$ -divergence ( $D_B^{(\beta)}$ ) via

$$\log_{\beta}(x) = \frac{1}{(\beta - 1)\beta} [\beta x^{\beta-1} - (\beta - 1)x^{\beta}]$$
$$D_B^{(\beta)}(g||p(\cdot|\theta)) = \mathbb{E}_g [\log_{\beta}(p(\mathbf{x}|\theta)) - \log_{\beta}(g(\mathbf{x}))]$$

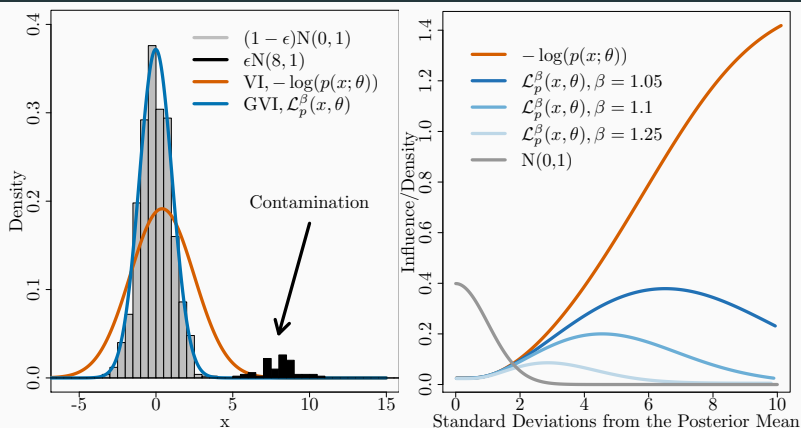
**Note 1:**  $D_B^{(\beta)} \rightarrow$  KLD as  $\beta \rightarrow 1!$

**Note 2:** Admits  $D_B^{(\beta)}$ -targeting loss as

$$\mathcal{L}_p^{\beta}(\theta, \mathbf{x}_i) = -\frac{1}{\beta - 1} p(\mathbf{x}_i|\theta)^{\beta-1} + \frac{l_{p,\beta}(\theta)}{\beta}, \quad l_{p,c}(\theta) = \int p(\mathbf{x}|\theta)^c d\mathbf{x}$$



## 4.2 GVI: The losses III/III



**Figure 5 – Left:** Robustness against model misspecification. Depicted are posterior predictives under  $\epsilon = 5\%$  outlier contamination using **VI** and  $P(\sum_{i=1}^n \mathcal{L}_p^\beta(\theta, x_i), \text{KLD}, \mathcal{Q})$ ,  $\beta = 1.5$ . **Right:** From Knoblauch et al. (2018). Influence of  $x_i$  on exact posteriors for different losses.

## 4.3 GVI: Uncertainty Quantification I/III

**GVI modularity: The uncertainty quantifier  $D$**

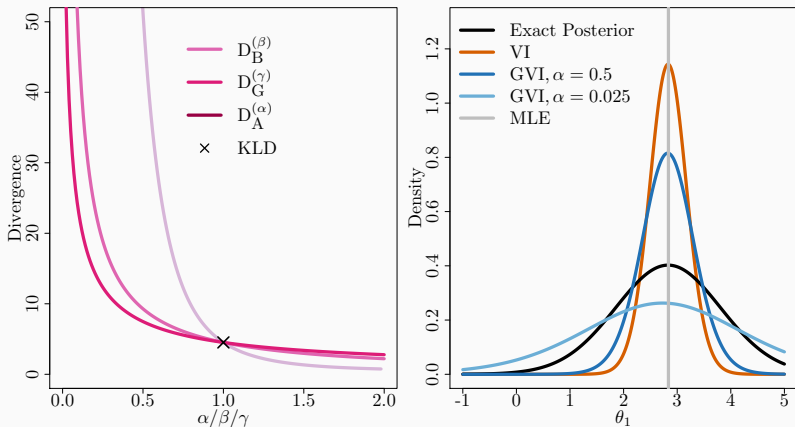
**Q:** Which **VI** drawbacks can be addressed via  $D$ ?

**A:** Any **uncertainty quantification** properties, e.g.

- Over-concentration ( = underestimating marginal variances)
- Sensitivity to badly specified priors
- ...

## 4.3 GVI: Uncertainty Quantification II/III

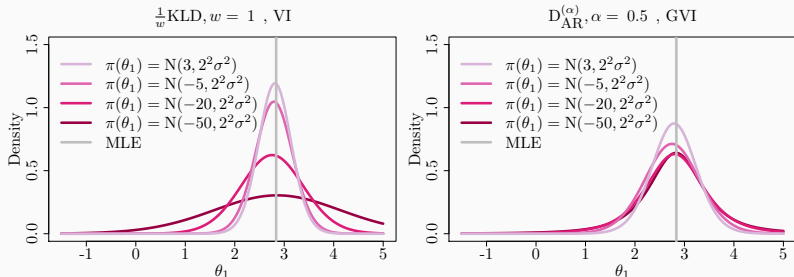
**Example 1:** GVI can fix over-concentrated posteriors



**Figure 6 – Left:** Magnitude of the penalty incurred by  $D(q||\pi)$  for different uncertainty quantifiers  $D$  and fixed densities  $\pi, q$ . **Right:** Using  $D_{AR}^{(\alpha)}$  with different choices of  $\alpha$  to “customize” uncertainty.

## 4.3 GVI: Uncertainty Quantification III/III

### Example 2: Avoiding prior sensitivity



**Figure 7** – Prior sensitivity with **VI** (left) vs. prior robustness with **GVI** (right). Priors are more badly specified for **dark**er shades.

## 4.4 GVI: Three use cases

**Summary:** GVI is natural in the M-open (i.e. real) world. Applications include

- (1) Robustness to model misspecification ( = adapting  $\ell_n$ )
- (2) “Customized” marginal variances ( = adapting  $D$ )
- (3) Prior robustness ( = adapting  $D$ )

### Purpose of part 5: GVI inference & experiments

- (1) How/when can we “black box” GVI?
- (2) A case study in robustness with Bayesian On-line Changepoint Detection
- (3) DVI vs GVI & changes in  $D$  (on Bayesian Neural Nets)
- (4) VI vs GVI & changes in  $\ell_n$  (on Deep Gaussian Processes)

## 5.1 Black Box GVI

**Setup:**  $\mathcal{Q} = \{q(\boldsymbol{\theta}|\boldsymbol{\kappa}) : \boldsymbol{\kappa} \in K\}$  variational family s.t.

- (i) one can sample  $\boldsymbol{\theta}^{(1:S)} \sim q(\boldsymbol{\theta}|\boldsymbol{\kappa})$ ;
- (ii) derivative  $\nabla_{\boldsymbol{\kappa}} \log(q(\boldsymbol{\theta}|\boldsymbol{\kappa}))$  exists.

**Case 1:** Closed form for  $\nabla_{\boldsymbol{\kappa}} D(q||\pi) \rightarrow$  unbiased estimate:

$$\nabla_{\boldsymbol{\kappa}} \hat{L}(q|\ell_n, D) = \frac{1}{S} \sum_{s=1}^S \left\{ \ell_n(\boldsymbol{\theta}^{(s)}, \mathbf{x}) \cdot \nabla_{\boldsymbol{\kappa}} \log(q(\boldsymbol{\theta}^{(s)}|\boldsymbol{\kappa})) \right\} + \nabla_{\boldsymbol{\kappa}} D(q||\pi)$$

**Thm. 7:** Closed forms for most  $\alpha/\beta/\gamma$ - & Rényi-divergence.

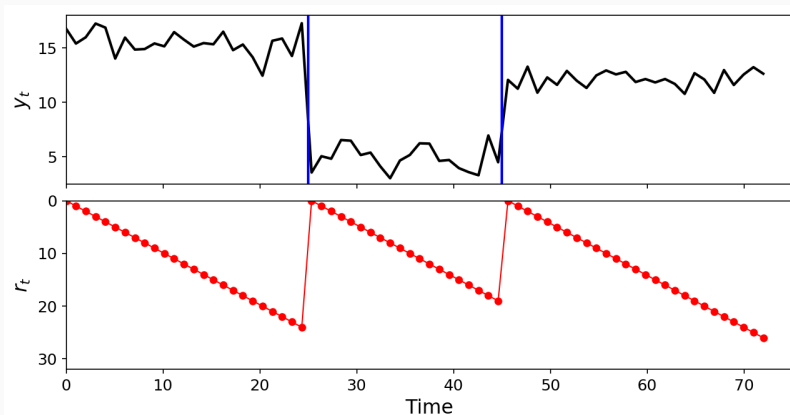
**Case 2:**  $D(q||\pi) = \mathbb{E}_q[\ell_{\boldsymbol{\kappa},\pi}^D(\boldsymbol{\theta})]$  (e.g.,  $f$ -divs)  $\rightarrow$  unbiased estimate:

$$\nabla_{\boldsymbol{\kappa}} \hat{L}(q|\ell_n, D) = \frac{1}{S} \sum_{s=1}^S \left\{ \left[ \ell_n(\boldsymbol{\theta}^{(s)}, \mathbf{x}) + \ell_{\boldsymbol{\kappa},\pi}^D(\boldsymbol{\theta}^{(s)}) \right] \cdot \nabla_{\boldsymbol{\kappa}} \log(q(\boldsymbol{\theta}^{(s)}|\boldsymbol{\kappa})) \right. \\ \left. + \nabla_{\boldsymbol{\kappa}} \ell_{\boldsymbol{\kappa},\pi}^D(\boldsymbol{\theta}^{(s)}) \right\}.$$

## 5.2 Standard BOCPD I/XI

Idea due to Adams and MacKay (2007) and Fearnhead and Liu (2007):

- (1) Define **Run-length at  $t = r_t$**   $\iff$  there was a CP at time  $t - r_t$ .
- (2) **Inference on last CP** via  $p(r_t|y_{1:t})$  rather than on *all* CPs
- (3) Resulting complexity:  $\mathcal{O}(t)$  **rather than**  $\mathcal{O}(\prod_{i=1}^t i)$ .





## 5.2 BOCPD + model selection (Knoblauch and Damoulas, 2018)

II/XI

**Idea:** Multiple models, different between segments

**New Random Variable:**  $m_t$ , the model at time  $t$

$$r_t | r_{t-1} \sim H(r_t, r_{t-1}) \quad [\text{conditional CP prior}] \quad (3a)$$

$$m_t | m_{t-1}, r_t \sim q(m_t | m_{t-1}, r_t) \quad [\text{conditional model prior}] \quad (3b)$$

$$\theta_m | m_t \sim \pi_{m_t}(\theta_{m_t}) \quad [\text{parameter prior}] \quad (3c)$$

$$y_t | m_t, \theta_{m_t} \sim f_{m_t}(y_t | \theta_{m_t}) \quad [\text{observation density}] \quad (3d)$$

where  $q(m_t | m_{t-1}, r_t) = \mathbb{1}_{\{r_t > 0\}} \delta(m_{t-1}) + \mathbb{1}_{\{r_t = 0\}} q(m_t)$ .

**Recursion:**

$$p(\mathbf{y}_1, r_1 = 0, m_1) = q(m_1) \int_{\Theta_{m_1}} f_{m_1}(\mathbf{y}_1 | \theta_{m_1}) \pi_{m_1}(\theta_{m_1}) d\theta_{m_1} = q(m_1) f_{m_1}(\mathbf{y}_1 | \mathbf{y}_0)$$

$$p(\mathbf{y}_{1:t}, r_t, m_t) = \sum_{m_{t-1}, r_{t-1}} \left\{ f_{m_t}(\mathbf{y}_t | \mathbf{y}_{1:(t-1)}, r_{t-1}) q(m_t | \mathbf{y}_{1:(t-1)}, r_{t-1}, m_{t-1}) \right. \\ \left. H(r_t, r_{t-1}) p(\mathbf{y}_{1:(t-1)}, r_{t-1}, m_{t-1}) \right\}$$

## 5.2 Standard BOCPD with model selection III/XI


$$p(\mathbf{y}_{1:t}, r_t, m_t) = \sum_{m_{t-1}, r_{t-1}} \left\{ f_{m_t}(\mathbf{y}_t | \mathbf{y}_{1:(t-1)}, r_{t-1}) q(m_t | \mathbf{y}_{1:(t-1)}, r_{t-1}, m_{t-1}) \right. \\ \left. H(r_t, r_{t-1}) p(\mathbf{y}_{1:(t-1)}, r_{t-1}, m_{t-1}) \right\}$$


### Inference:

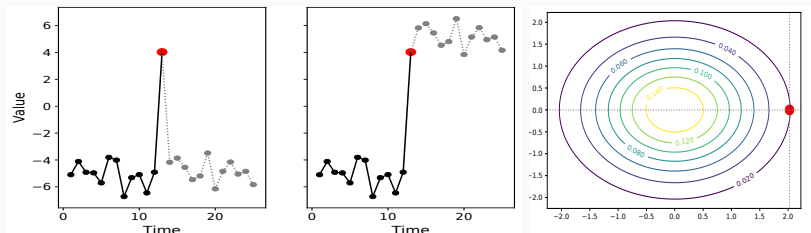
- (1) Evidence:  $p(\mathbf{y}_{1:t}) = \sum_{r_t, m_t} p(\mathbf{y}_{1:t}, r_t, m_t)$
- (2) run-length & model posterior:  $p(r_t, m_t | \mathbf{y}_{1:t}) = p(\mathbf{y}_{1:t}, r_t, m_t) / p(\mathbf{y}_{1:t})$
- (3) Prediction:  $p(\mathbf{y}_{t+1} | \mathbf{y}_{1:t}) = \sum_{r_t, m_t} f_{m_t}(\mathbf{y}_{t+1} | \mathbf{y}_{1:t}, r_t) p(r_t, m_t | \mathbf{y}_{1:t})$
- (4) Run-length marginal posterior:  $p(r_t | \mathbf{y}_{1:t}) = \sum_{m_t} p(r_t, m_t | \mathbf{y}_{1:t})$
- (5) Model marginal posterior:  $p(m_t | \mathbf{y}_{1:t}) = \sum_{r_t} p(r_t, m_t | \mathbf{y}_{1:t})$ .
- (6) MAP segmentation:  
 $MAP_t = \max_{r,t} \{ MAP_{t-r-1} \cdot p(r_t = r, m_t = m | \mathbf{y}_{1:t}) \}$

## 5.2 Issue: Outliers & misspecification IV/X

Last plot: Clear that model selection non-robust! Why?

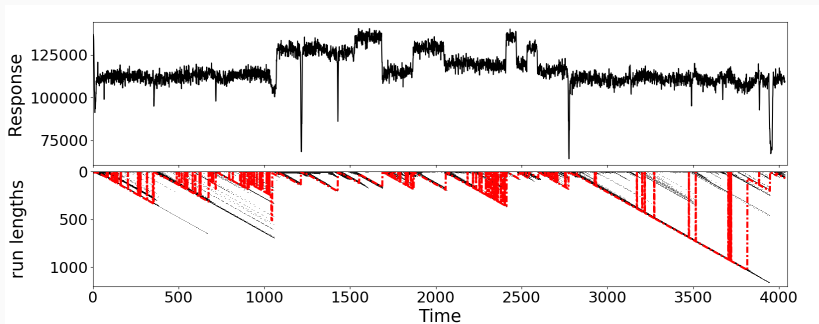
 On-line processing

 Moderate/high dimensions for  $y_t$



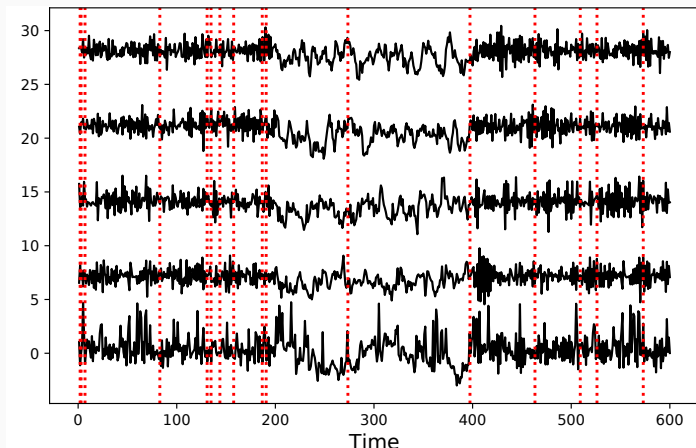
**Figure 8 – Left, Center:** Price for on-line processing is that outliers are confused with changepoints. **Right:** Multivariate densities become very small even if outliers occur only in a single dimension.

## 5.2 Issue: Outliers & misspecification V/X



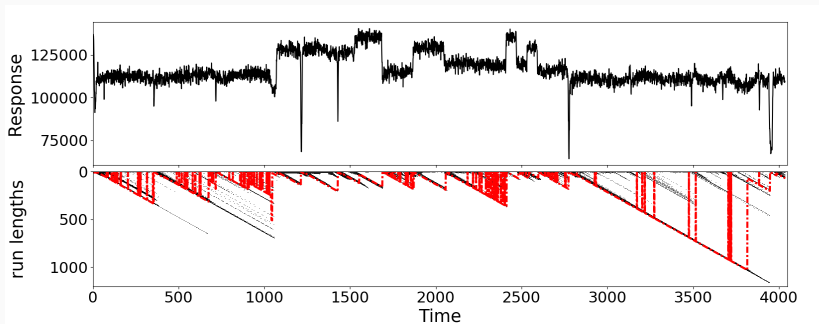
## 5.2 Issue: Outliers & misspecification VI/X

Five Autoregressive processes with two CPs

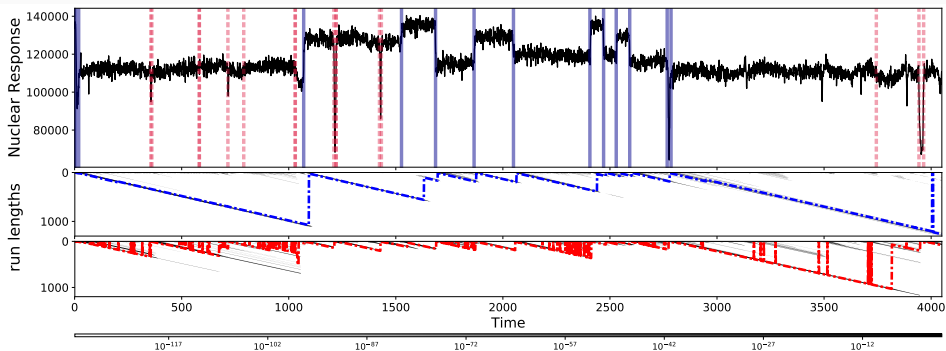


**Figure 9** – Maximum A Posteriori (MAP) CPs of **standard** BOCPD shown as dashed vertical lines. True CPs at  $t = 200, 400$ .

## 5.2 Fix: Adapt loss (Knoblauch et al., 2018) VII/XI



## 5.2 Fix: Robustness by adapting the loss VIII/XI

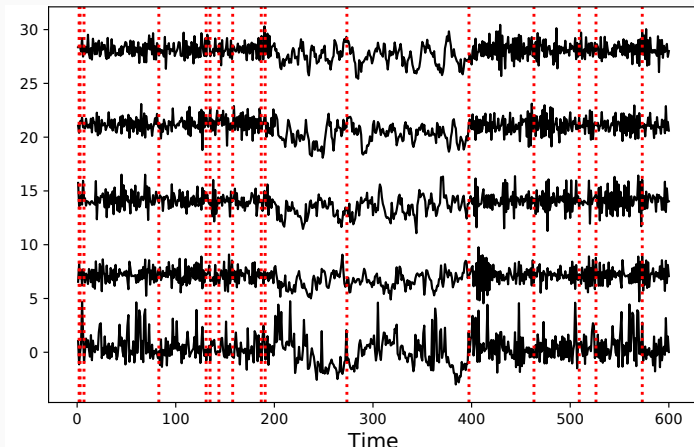


**Figure 10 – Robust segmentation and run-length distribution and additionally found CPs with non-robust run-length distribution**

[FDR:  $> 99\% \implies 8\%$  and reduction in MSE (MAE) by 10% (6%)]

## 5.2 Fix: Robustness by adapting the loss IX/XII

Five Autoregressive processes with two CPs

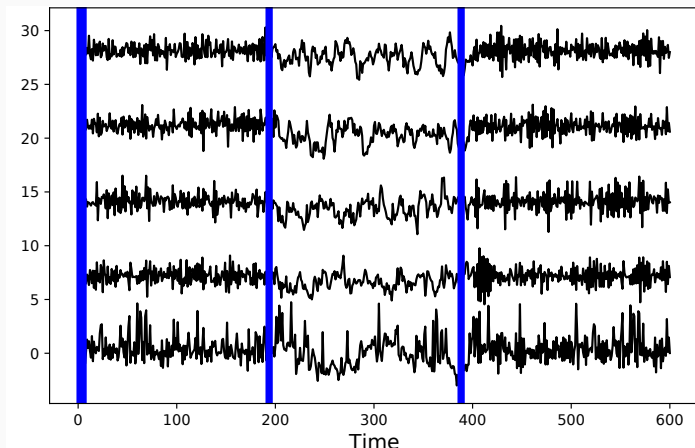


**Figure 11** – Maximum A Posteriori (MAP) CPs of **standard** BOCPD shown as dashed vertical lines. True CPs at  $t = 200, 400$ .



## 5.2 Fix: Robustness by adapting the loss $X/XI$

Five Autoregressive processes with two CPs



**Figure 12** – Maximum A Posteriori (MAP) CPs of **robust** BOCPD shown as solid vertical lines. True CPs at  $t = 200, 400$ .

## 5.2 Fix: Robustness by adapting the loss $XI/XI$

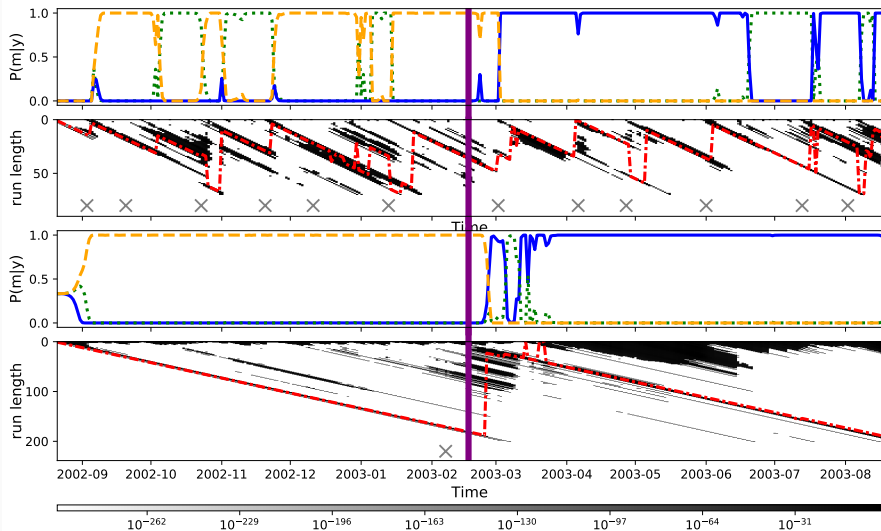


Figure 13 – Top & bottom two panels: standard & robust BOCPD.

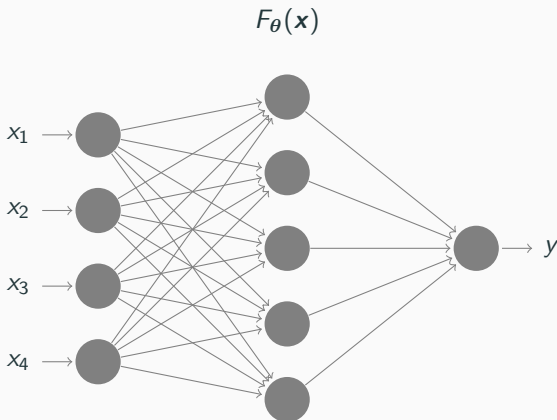
## 5.3 Experiments with Bayesian Neural Nets (BNNs) I/IV

BNNs are intractable Bayesian regression models with

$$y|\mathbf{x} \sim \mathcal{N}(\mathbf{y}; F_{\theta}(\mathbf{x}), \sigma^2),$$

with  $F_{\theta}(\mathbf{x})$  defining a non-linear transform of  $\mathbf{x}$  parameterized by  $\theta$ .

(**Note:** Our experiments use one hidden layer with 50 ReLU neurons.)



## 5.3 Experiments with Bayesian Neural Nets (BNNs) II/IV

**Methods:** Comparison of black box approximate Bayesian methods:

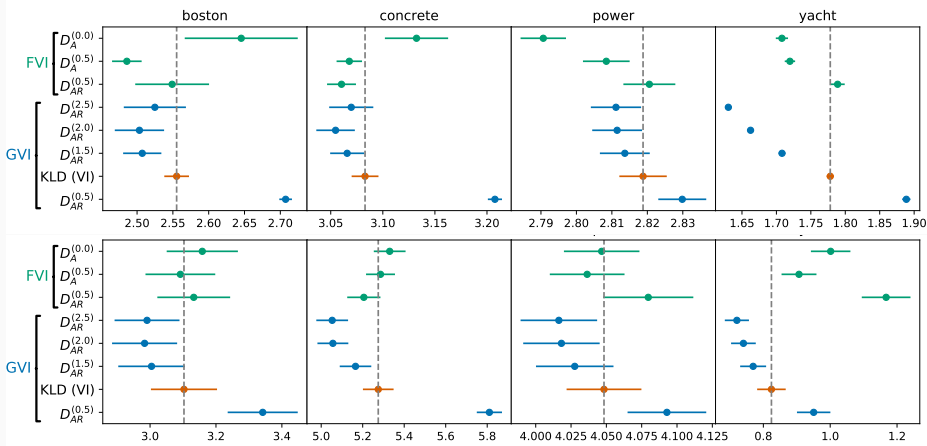
- **VI**
- **DVI** based on  $F = D_{AR}^{(\alpha)}$  (Li and Turner, 2016)
- **DVI** based on  $F = D_A^{(\alpha)}$  (Hernández-Lobato et al., 2016)
- **GVI** with  $D = D_{AR}^{(\alpha)}$ .

**Note: Everything run with settings of Li and Turner (2016) and Hernández-Lobato et al. (2016)**

- Variational family  $\mathcal{Q}$ : A fully factorized normal
- Optimization of  $\sigma^2$  (i.e., point estimation akin to type-II ML)
- ADAM (Kingma and Ba, 2014) with default settings and 500 epochs
- 50 Random splits with 90:10 training:test ratio
- benchmark UCI (Lichman et al., 2013) datasets

...

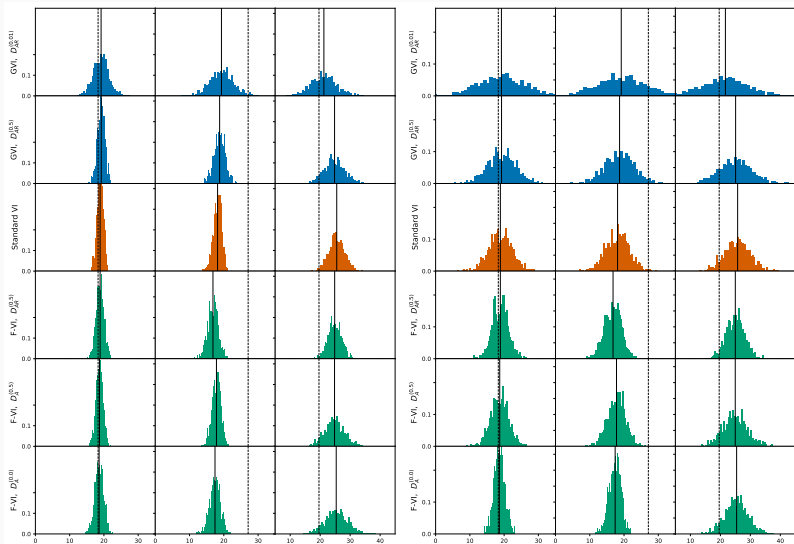
## 5.3 Experiments with Bayesian Neural Nets (BNNs) III/IV



**Figure 14** – Performance on BNNs: **DVI**, **GVI** with  $D = D_{AR}^{(\alpha)}$ , and **VI**. **Top**: Negative test log likelihoods. **Bottom row**: Test RMSE.

**Observation**: **GVI** outperforms **VI** for **over-concentrated** posteriors (i.e.  $\alpha > 1$ )! So how does **under-concentrated DVI** outperform **VI**?!?

## 5.3 Experiments with Bayesian Neural Nets (BNNs) IV/V



**Figure 15** – Left: Parameter posteriors (DVI as expected). Right: Posterior predictives (DVI not as expected)

## 5.3 Experiments with Bayesian Neural Nets (BNNs) V/V

**Q:** Why does this happen for **DVI** and not for **GVI**?!

**A:** **DVI** does not distinguish uncertainty quantification & loss!

**DVI objective:**  $\sigma^2$  affects **target** (!)

$$\hat{\sigma}^2, q^*(\theta|\hat{\sigma}^2, \kappa) = \arg \min_{\sigma^2} \left\{ \arg \min_{q \in \mathcal{Q}} F \left( q(\theta|\sigma^2, \kappa) \parallel \underbrace{\tilde{q}(\theta|\sigma^2, \mathbf{x}, \mathbf{y})}_{\text{i.e., } \tilde{q}=\tilde{q}^\sigma} \right) \right\}$$

$\Rightarrow$  optimizing for  $\sigma^2 =$  **changing the target**  $\tilde{q}^\sigma$ !

**GVI objective:**  $\sigma^2$  indexes the **loss** only

$$\hat{\sigma}^2, q^*(\theta|\hat{\sigma}^2, \mathbf{x}, \mathbf{y}) = \arg \min_{\sigma^2} \left\{ \arg \min_{q \in \mathcal{Q}} \left\{ \mathbb{E}_q \left[ \underbrace{\ell_n(\theta, \mathbf{x}|\mathbf{y}, \sigma^2)}_{\text{i.e., } \ell_n=\ell_n^\sigma} \right] + D(q|\pi) \right\} \right\}$$

$\Rightarrow$  optimizing for  $\sigma^2 =$  finding **optimal loss**  $\ell_n^\sigma$

## 5.4 Experiments with Deep Gaussian Processes (DGPs) I/II

**Principal idea:** Use the BNN architecture with GP priors on  $F_\theta(\cdot)$ :

$$\begin{aligned}y|\mathbf{F}^L &\sim p(y|\mathbf{F}^L) \\ \mathbf{F}^L|\mathbf{F}^{L-1} &\sim \text{GP}\left(\mu^L(\mathbf{F}^{L-1}), \mathbf{K}^L(\mathbf{F}^{L-1}, \mathbf{F}^{L-1})\right) \\ \mathbf{F}^{L-1}|\mathbf{F}^{L-2} &\sim \text{GP}\left(\mu^{L-1}(\mathbf{F}^{L-2}), \mathbf{K}^{L-1}(\mathbf{F}^{L-2}, \mathbf{F}^{L-2})\right) \\ &\dots \\ \mathbf{F}^1|\mathbf{x} &\sim \text{GP}\left(\mu^1(\mathbf{x}), \mathbf{K}^1(\mathbf{x}, \mathbf{x})\right),\end{aligned}$$

**Methods:** Comparison of black box approximate Bayesian methods:

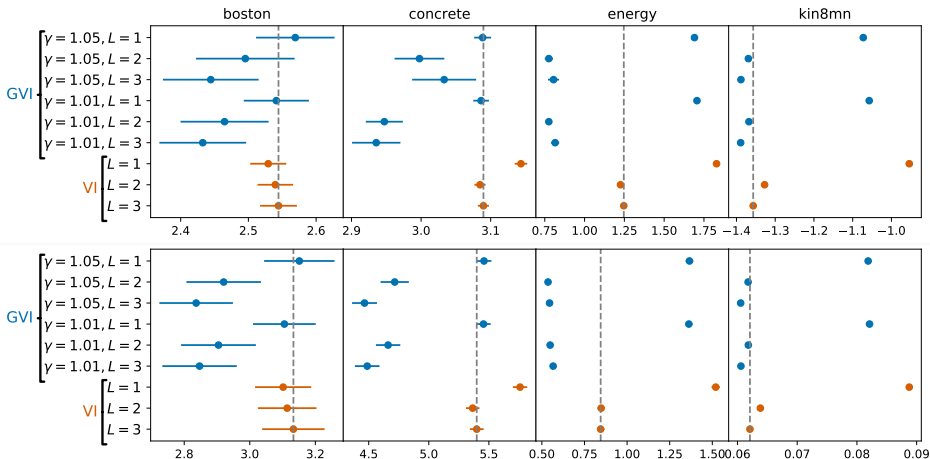
- State of the art **VI** (Salimbeni and Deisenroth, 2017)  
(comprehensively beat competing **DVI** methods (Bui et al., 2016))
- **GVI** with  $\ell_n = \sum_{i=1}^n \mathcal{L}_p^\gamma(\boldsymbol{\theta}, \mathbf{x}_i)$ .

**Note: Everything run with settings of Salimbeni and Deisenroth (2017)**

[Derivations for **DGP-GVI**: <https://arxiv.org/abs/1904.02303>.]



## 5.4 Experiments with Deep Gaussian Processes (DGPs) II/II



**Figure 16** – DGP performance with  $L$  layers, **GVI** with & **VI**. **Top row**: Negative test log likelihoods. **Bottom row**: Test RMSE.

# Summary & Conclusion

## Summary:

- Part 1: Ways to look at Bayesian inference: belief updates (about arbitrary parameters) & optimization over  $\mathcal{P}(\Theta)$
- Part 2: Bayesian inference as a modular & interpretable triplet  $P(\ell_n, D, \Pi)$ : **loss**, **uncertainty quantifier** & **admissible posteriors**.
- Part 3: Fallout of  $P(\ell_n, D, \Pi)$ : **VI** optimality & **DVI** suboptimality  $\rightarrow$  **GVI**
- Part 4: Some of **GVI**'s use cases: Robust losses, alternative ways of quantifying uncertainty. Also: its upper bound interpretation
- Part 5: Black box methods with **GVI** & empirical performance.

## Main Conclusions:

- (I) **GVI**: **principled & explicit** design of  $Q$ -constrained posteriors
- (II) **GVI**: **tackles drawbacks of VI** (e.g., robustness, marginals)
- (III) **GVI**: **State of the art**  $Q$ -constrained posteriors **on BNNs & DGPs**

# Main References i

- Adams, R. P. and MacKay, D. J. (2007). Bayesian online changepoint detection. *arXiv preprint arXiv:0710.3742*.
- Alquier, P. and Guedj, B. (2018). Simpler pac-bayesian bounds for hostile data. *Machine Learning*, 107(5):887–902.
- Alquier, P., Ridgway, J., and Chopin, N. (2016). On the properties of variational approximations of gibbs posteriors. *The Journal of Machine Learning Research*, 17(1):8374–8414.
- Ambrogioni, L., Güçlü, U., Güçlütürk, Y., Hinne, M., van Gerven, M. A. J., and Maris, E. (2018). Wasserstein variational inference. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 31*, pages 2478–2487. Curran Associates, Inc.
- Bissiri, P. G., Holmes, C. C., and Walker, S. G. (2016). A general framework for updating belief distributions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):1103–1130.
- Bui, T., Hernández-Lobato, D., Hernandez-Lobato, J., Li, Y., and Turner, R. (2016). Deep gaussian processes for regression using approximate expectation propagation. In *International Conference on Machine Learning*, pages 1472–1481.
- Dieng, A. B., Tran, D., Ranganath, R., Paisley, J., and Blei, D. (2017). Variational inference via  $\chi$  upper bound minimization. In *Advances in Neural Information Processing Systems*, pages 2732–2741.
- Fearnhead, P. and Liu, Z. (2007). On-line inference for multiple changepoint problems. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(4):589–605.
- Futami, F., Sato, I., and Sugiyama, M. (2017). Variational inference based on robust divergences. *arXiv preprint arXiv:1710.06595*.
- Ganchev, K., Gillenwater, J., Taskar, B., et al. (2010). Posterior regularization for structured latent variable models. *Journal of Machine Learning Research*, 11(Jul):2001–2049.
- Ghosh, A. and Basu, A. (2016). Robust bayes estimation using the density power divergence. *Annals of the Institute of Statistical Mathematics*, 68(2):413–437.
- Grünwald, P., Van Ommen, T., et al. (2017). Inconsistency of bayesian inference for misspecified linear models, and a proposal for repairing it. *Bayesian Analysis*, 12(4):1069–1103.
- Hernández-Lobato, J. M., Li, Y., Rowland, M., Hernández-Lobato, D., Bui, T. D., and Turner, R. E. (2016). Black-box  $\alpha$ -divergence minimization. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning-Volume 48*, pages 1511–1520.
- Holmes, C. and Walker, S. (2017). Assigning a value to a power likelihood in a general bayesian model. *Biometrika*, 104(2):497–503.
- Hooker, G. and Vidyashankar, A. N. (2014). Bayesian model robustness via disparities. *Test*, 23(3):556–584.
- Huang, C.-W., Tan, S., Lacoste, A., and Courville, A. C. (2018). Improving explorability in variational inference with annealed variational objectives. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 31*, pages 9724–9734. Curran Associates, Inc.

# Main References ii

- Jewson, J., Smith, J., and Holmes, C. (2018). Principles of Bayesian inference using general divergence criteria. *Entropy*, 20(6):442.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Knoblauch, J. and Damoulas, T. (2018). Spatio-temporal Bayesian on-line changepoint detection with model selection. In *Proceedings of the 27th International Conference on Machine Learning (ICML-18)*.
- Knoblauch, J., Jewson, J., and Damoulas, T. (2018). Doubly robust Bayesian inference for non-stationary streaming data using  $\beta$ -divergences. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 64–75.
- Li, Y. and Turner, R. E. (2016). Rényi divergence variational inference. In *Advances in Neural Information Processing Systems*, pages 1073–1081.
- Lichman, M. et al. (2013). Uci machine learning repository.
- Miller, J. W. and Dunson, D. B. (2018). Robust bayesian inference via coarsening. *Journal of the American Statistical Association*, (just-accepted):1–31.
- Minka, T. P. (2001). Expectation propagation for approximate bayesian inference. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, pages 362–369. Morgan Kaufmann Publishers Inc.
- Opper, M. and Winther, O. (2000). Gaussian processes for classification: Mean-field algorithms. *Neural computation*, 12(11):2655–2684.
- Ranganath, R., Tran, D., Altsosaar, J., and Blei, D. (2016). Operator variational inference. In *Advances in Neural Information Processing Systems*, pages 496–504.
- Regli, J.-B. and Silva, R. (2018). Alpha-beta divergence for variational inference. *arXiv preprint arXiv:1805.01045*.
- Saha, A., Bharath, K., and Kurtek, S. (2017). A geometric variational approach to bayesian inference. *arXiv preprint arXiv:1707.09714*.
- Salimbeni, H. and Deisenroth, M. (2017). Doubly stochastic variational inference for deep gaussian processes. In *Advances in Neural Information Processing Systems*, pages 4588–4599.
- Yang, Y., Pati, D., and Bhattacharya, A. (2017).  $\alpha$ -variational inference with statistical guarantees. *arXiv preprint arXiv:1710.03266*.
- Zellner, A. (1988). Optimal information processing and bayes's theorem. *The American Statistician*, 42(4):278–280.
- Zhu, J., Chen, N., and Xing, E. P. (2014). Bayesian inference with posterior regularization and applications to infinite latent svms. *The Journal of Machine Learning Research*, 15(1):1799–1847.

## 4.2 GVI: Choosing loss hyperparameters

$$\begin{aligned}\mathcal{L}_p^\beta(\boldsymbol{\theta}, x_i) &= -\frac{1}{\beta-1} p(x_i|\boldsymbol{\theta})^{\beta-1} + \frac{I_{p,\beta}(\boldsymbol{\theta})}{\beta} \\ \mathcal{L}_p^\gamma(\boldsymbol{\theta}, x_i) &= -\frac{1}{\gamma-1} p(x_i|\boldsymbol{\theta})^{\gamma-1} \frac{\gamma}{I_{p,\gamma}(\boldsymbol{\theta})^{\frac{\gamma-1}{\gamma}}} \\ I_{p,c}(\boldsymbol{\theta}) &= \int p(x|\boldsymbol{\theta})^c dx\end{aligned}$$

where  $I_{p,c}(\boldsymbol{\theta}) = \int p(x|\boldsymbol{\theta})^c dx$ .

**Note 1:**  $\mathcal{L}_p^\gamma(\boldsymbol{\theta}, x_i)$  multiplicative & always  $< 0 \rightarrow$  store as log!

**Note 2:** Conditional independence  $\neq$  additive for  $\mathcal{L}_p^\beta(\boldsymbol{\theta}, x_i), \mathcal{L}_p^\gamma(\boldsymbol{\theta}, x_i)$

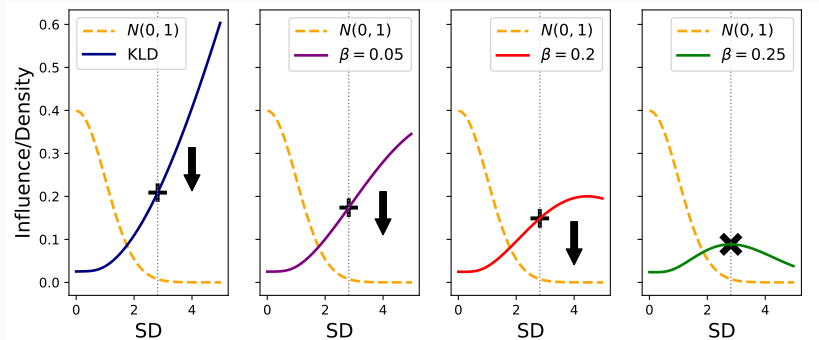
**Note 3:** In practice, usually best to choose  $\beta/\gamma = 1 + \varepsilon$  for some small  $\varepsilon$

## Appendix: Choosing hyperparameters

**Q:** Any principled way of choosing hyperparameters?

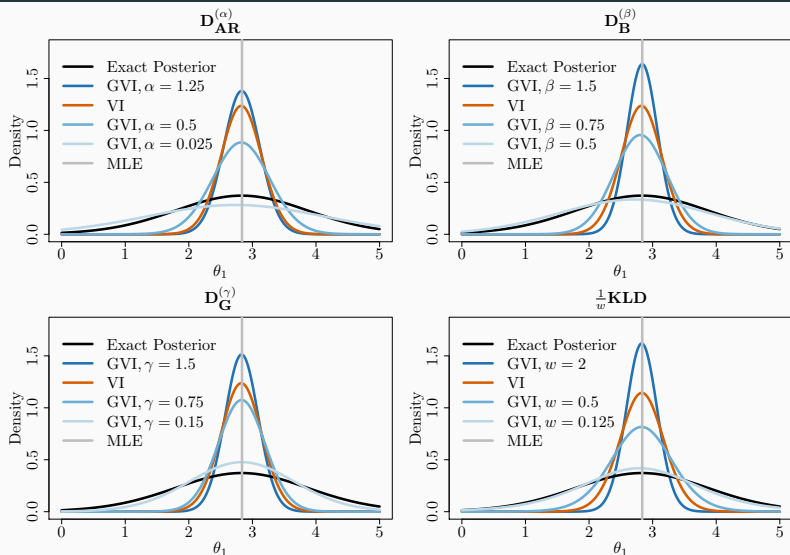
**A:** Very much unsolved problem, solutions so far:

- $D$ : brute force (CV) (Regli and Silva, 2018) [slow/expensive]
- $\ell_n$ : Via *points of highest influence* (Knoblauch et al., 2018)
- $\ell_n$ : on-line updates using loss-minimization (Knoblauch et al., 2018)



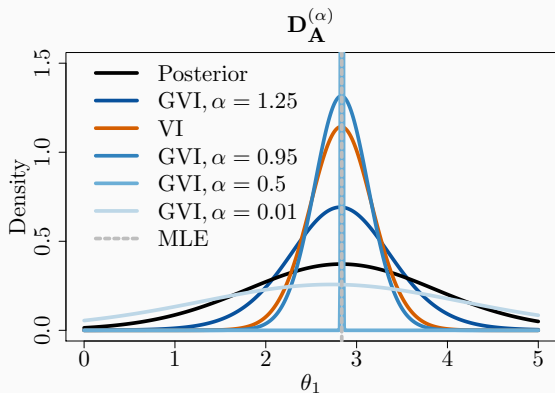
**Figure 17** – Illustration of the initialization procedure using *points of highest influence* logic, from left to right.

# Appendix: Choosing $D$ for conservative marginals I/II



**Figure 18** – Marginal **vi** and **GVI** posterior for a Bayesian linear model under the  $D_{AR}^{(\alpha)}$ ,  $D_B^{(\beta)}$ ,  $D_G^{(\gamma)}$  and  $\frac{1}{w} \text{KLD}$  uncertainty quantifier for different values of the divergence hyperparameters.

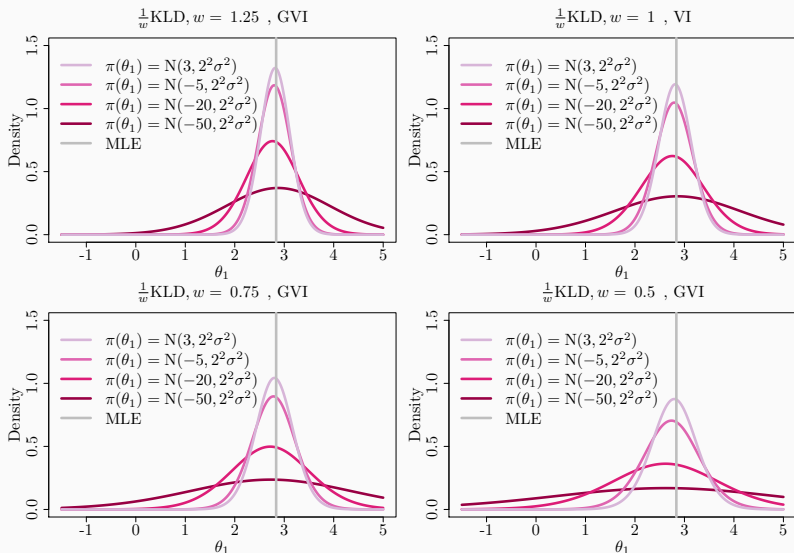
## Appendix: Choosing $D$ for conservative marginals II/II



**Figure 19** – Marginal **VI** and **GVI** posterior for a Bayesian linear model under the  $D_A^{(\alpha)}$  uncertainty quantifier. The boundedness of the  $D_A^{(\alpha)}$  causes **GVI** to severely over-concentrate if  $\alpha$  is not carefully specified.

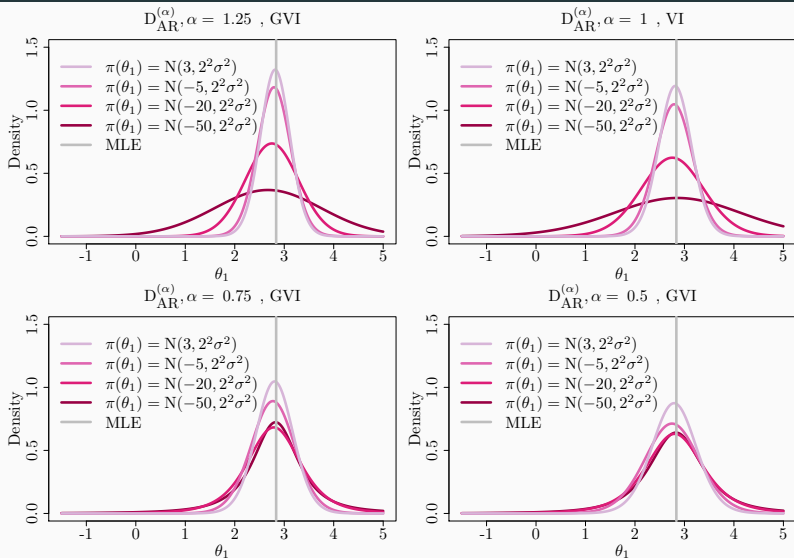


# Appendix: Choosing $D$ for prior robustness I/IV



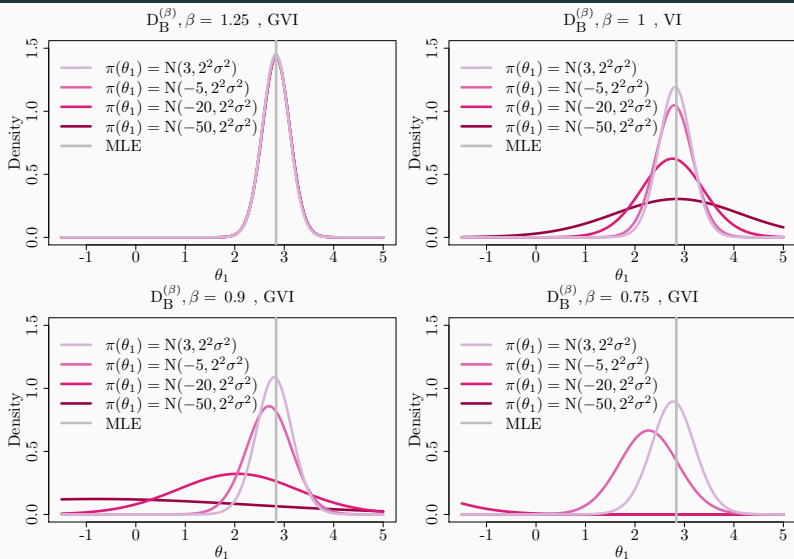
**Figure 20** – Marginal **VI** and **GVI** posterior for a Bayesian linear model under different priors, using  $D = \frac{1}{w}\text{KLD}$  as the uncertainty quantifier.

# Appendix: Choosing $D$ for prior robustness II/IV



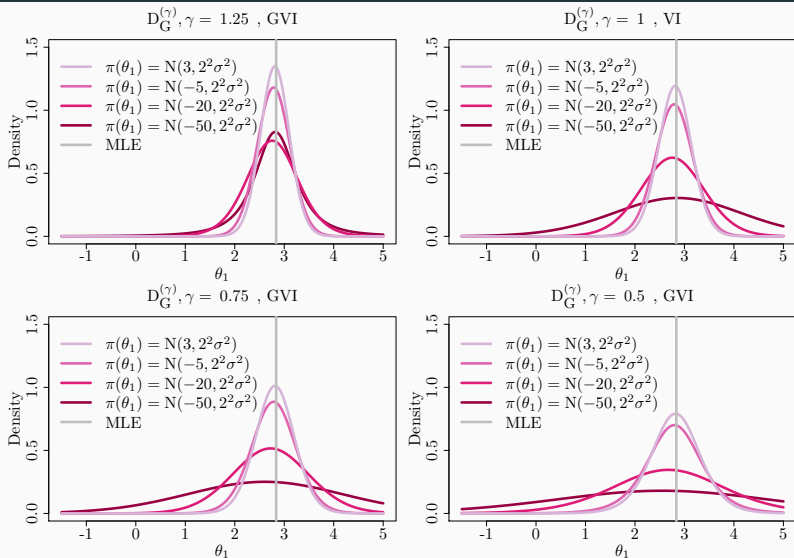
**Figure 21** – Marginal **VI** and **GVI** posterior for a Bayesian linear model under different priors, using  $D = D_{AR}^{(\alpha)}$  as the uncertainty quantifier.

# Appendix: Choosing $D$ for prior robustness III/IV



**Figure 22** – Marginal **VI** and **GVI** posterior for a Bayesian linear model under different priors, using  $D = D_B^{(\beta)}$  as the uncertainty quantifier.

# Appendix: Choosing $D$ for prior robustness IV/IV



**Figure 23** – Marginal **VI** and **GVI** posterior for a Bayesian linear model under different priors, using  $D = D_G^{(\gamma)}$  as the uncertainty quantifier.

## Appendix: GVI lower bound interpretation I/II

**Question:** VI is also interpretable as optimizing a lower bound on the evidence! Is there anything comparable for GVI?

**Answer:** Yes, e.g. for  $D_B^{(\beta)}$ ,  $D_G^{(\gamma)}$ ,  $D_{AR}^{(\alpha)}$ : Consider *generalized* evidence:

**Recall:** Generalized Bayes posterior (Bissiri et al., 2016) is

$$q_{\ell_n}^*(\theta) \propto \pi(\theta) \exp\{-\ell_n(\theta, \mathbf{x})\} \quad \text{and so } p_{\ell_n}(\mathbf{x}) = \int_{\Theta} q_{\ell_n}^*(\theta) d\theta$$

GVI's objectives  $L(q|\mathbf{x}, D, \ell_n)$  will optimize

$$L(q|\mathbf{x}, D, \ell_n) \geq g^D\left(\underbrace{-\log p_{f^D(\ell_n)}(\mathbf{x})}_{\substack{\text{negative log evidence;} \\ f^D(\ell_n) \text{ maps } \ell_n \text{ into a new loss}}}\right) + \underbrace{T^D(q)}_{\text{Approximate target}}$$

(**Note:** VI is special case where this holds with *equality* (so that the approximate target is the exact target) and where  $g^{\text{KLD}}(\mathbf{x}) = \mathbf{x}$ ,  $L(q|\mathbf{x}, D, \ell_n) = \text{ELBO}(q)$ ,  $T^{\text{KLD}}(q) = \text{KLD}(q||q_{\ell_n}^*)$ ,  $f^{\text{KLD}}(\ell_n) = \ell_n$ .)

## Appendix: GVI lower bound interpretation II/II

GVI's objectives  $L(q|\mathbf{x}, D, \ell_n)$  will optimize

$$L(q|\mathbf{x}, D, \ell_n) \geq g^D\left(\underbrace{-\log p_{f^D(\ell_n)}(\mathbf{x})}_{\substack{\text{negative log evidence;} \\ f^D(\ell_n) \text{ maps } \ell_n \text{ into a new loss}}}\right) + \underbrace{T^D(q)}_{\text{Approximate target}}$$

**Example:** Rényi's  $\alpha$ -divergence ( $D_{AR}^{(\alpha)}$ ) for  $\alpha > 1$  gives

$$g^{D_{AR}^{(\alpha)}}(\mathbf{x}) = \frac{1}{\alpha} \mathbf{x},$$

$$f^{D_{AR}^{(\alpha)}}(\ell_n) = \alpha \ell_n,$$

$$T_{D_{AR}^{(\alpha)}}(q) = \frac{1}{\alpha} \text{KLD}(q||q_{\alpha \ell_n}^*),$$

so putting it together one finds that for  $D = D_{AR}^{(\alpha)}$  with  $\alpha > 1$ ,

$$L(q|\mathbf{x}, D, \ell_n) \geq -\frac{1}{\alpha} \log p_{\alpha \ell_n}(\mathbf{x}) + \frac{1}{\alpha} \text{KLD}(q||q_{\alpha \ell_n}^*)$$

(Which is just a  $\frac{1}{\alpha}$ -scaled version of the ELBO for the loss  $\alpha \ell_n$ !)

# 5.4 Experiments with Deep Gaussian Processes (DGPs) III/IV

