



Optimization-centric Generalizations of Bayesian Inference

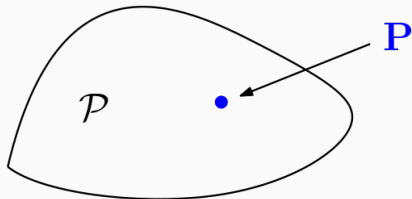
Generalized Variational Inference & beyond

November 5, 2021

Biometrika Fellow @ Department of Statistical Science, UCL
Visiting Researcher @ The Alan Turing Institute for Data Science and AI

What you should take away from today

- (1) Bayesian inference = a **single** optimization problem $P \in \mathcal{P}$
- (2) Optimization problem $P \in \mathcal{P}$ needs **strict assumptions**
- (3) These assumptions are often violated in Machine Learning (ML)
- (4) **Many** other optimization problems $P' \in \mathcal{P}$ you could solve instead
- (5) In ML, we would often do better to solve $P' \in \mathcal{P}$ rather than P



- (1) **All about that Bayes**
 - (1.1) Bayes' rule
 - (1.2) Assumptions of Bayesian inference
 - (1.3) Machine Learning, Assumptions & Bayesian inference

- (2) **An optimization-centric generalization**
 - (2.1) Bayesian inference is optimization
 - (2.2) Modularity & Interpretation
 - (2.3) A natural generalization
 - (2.4) Special cases of note
 - (2.5) Generalized Variational Inference

- (3) **Properties of this generalization**
 - (3.1) 'Sanity checks': existence, uniqueness, consistency
 - (3.2) Axiomatic foundations
 - (3.3) Relation to VI, PAC-Bayes bounds, power posteriors, ...
 - (3.4) Closed Forms

- (4) **Applications**
 - (4.1) Robustness: prior misspecification
 - (4.2) Robustness: likelihood misspecification
 - (4.3) Simplifying Bayesian computation

Literature I will talk about

Main focus (parts 1+2):

JK, Jack Jewson, & Theo Damoulas; *Generalized Variational Inference: Three Arguments for deriving new Posteriors*, currently minor revisions at JMLR <https://arxiv.org/abs/1904.02063>

literature I touch upon (parts 3+4):

Takuo Matsubara, JK, Francois-Xavier Briol, & Chris Oates; *Generalised Bayesian Inference with Stein Discrepancies: Robust Bayes for Models with an Intractable Likelihood*, submitted to JRSS-B 2021 <https://arxiv.org/abs/2104.07359>

JK; *Frequentist Consistency of Generalized Variational Inference*, 2019 <https://arxiv.org/abs/1904.04946>

JK; *Robust Deep Gaussian Processes*, 2019 <https://arxiv.org/abs/1904.02303>

JK, Jack Jewson, & Theo Damoulas; *Doubly Robust Bayesian Inference for Non-Stationary Streaming Data using β -Divergences*, NeurIPS 2018 <https://arxiv.org/abs/1806.02261>

Sebastian Schmon, Patrick Cannon, & JK; *Generalized Posteriors in Approximate Bayesian Computation*, AABI 2021 <https://arxiv.org/abs/2011.08644>

Pierre Alquier; *Non-exponentially weighted aggregation: regret bounds for unbounded losses*, 2020 <https://arxiv.org/abs/2009.03017>

Non-exhaustive list of important literature I do not touch upon:

Alexander A. Alemi; *Variational Predictive Information Bottleneck*; AABI 2019 <https://arxiv.org/abs/1910.10831>

Jeffrey Miller; *Asymptotic normality, concentration, and coverage of generalized posteriors*, working paper 2019 (arxiv 1907.09611)

Badr-Eddine Chérif-Abdellatif; *Contributions to the theoretical study of variational inference and robustness*, PhD thesis 2020

...

(1) All about that Bayes

Not actually Rev. Thomas Bayes



Not actually the form of the result in Bayes (1763), but in Laplace (1774)

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

(1.1) Bayes' Rule: Interpretation as belief updates

Ingredients:

- n observations $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$,
- **prior** $\pi(\theta)$,
- **likelihood terms** $\{p(x_i|\theta)\}_{i=1}^n$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Output (via Bayes' Rule) = **posterior belief**:

$$\underbrace{q_n^*(\theta)}_{\text{posterior}} \propto \underbrace{\pi(\theta)}_{\text{prior}} \prod_{i=1}^n \underbrace{p(x_i|\theta)}_{\text{likelihood terms}} \quad (1)$$

Inference interpretation = **belief updates**:

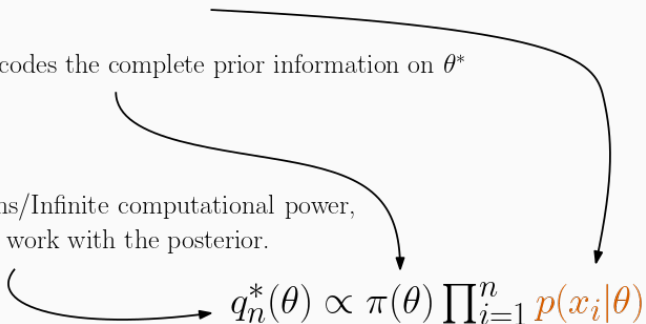
$$\pi(\theta) \xrightarrow{\text{Update with Bayes' rule via } \{p(x_i|\theta)\}_{i=1}^n} q_n^*(\theta)$$

(1.2) Three assumptions underlying Bayesian inference

1.) Likelihoods are correct, i.e. $p(x_i|\theta^*) \stackrel{!}{=} d\mathbb{P}(x_i)$ for some unknown θ^*

2.) Prior encodes the complete prior information on θ^*

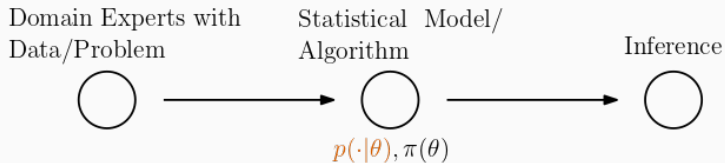
3.) Closed forms/Infinite computational power, so that one can work with the posterior.



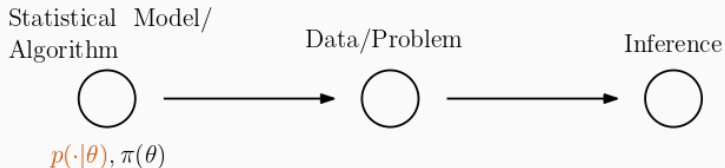
$q_n^*(\theta) \propto \pi(\theta) \prod_{i=1}^n p(x_i|\theta)$

(1.3) Machine Learning & Bayesian inference

Traditional role of Statistics in science:



Modern Machine Learning:



(1.3) Machine Learning & Bayesian inference

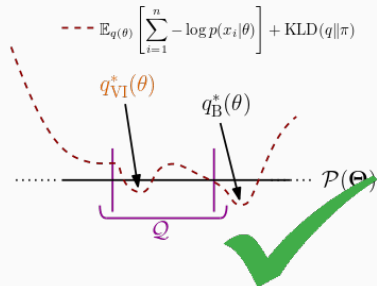
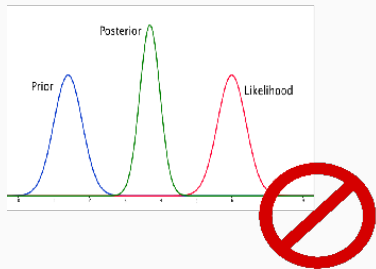
Conclusion I: Assumptions of Bayesian inference fit traditional science

- (1) **Correct Likelihoods:** Domain Experts help design these
- (2) **Priors encoding complete information:** Experiments/studies build on **prior** work, which can often be summarized in distributions.
(Prior elicitation is an entire field within Bayesian statistics.)
- (3) **Closed forms/infinite computational power:** Data collection is expensive; computational cost much less important

Conclusion II: ML violates assumptions of Bayesian Inference

- (1) **Correct Likelihoods:** Likelihoods are defined before data is seen
- (2) **Priors encoding complete information:** Priors often impossible to elicit (e.g., Bayesian Neural Networks)
- (3) **Closed forms/infinite computational power:** Violated by virtually any ML application

(2) An optimization-centric generalization



(2.1) Bayesian inference $\stackrel{!}{=}$ optimization

The Bayes posterior $q_n^*(\theta) \propto \pi(\theta) \prod_{i=1}^n p(x_i|\theta)$ uniquely solves

$$q_n^*(\theta) = \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \underbrace{\mathbb{E}_{q(\theta)} \left[\sum_{i=1}^n -\log(p(x_i|\theta)) \right]}_{\text{minimized by } q(\theta) = \delta_{\hat{\theta}_n}(\theta), \hat{\theta}_n = \text{MLE}} + \underbrace{\text{KLD}(q||\pi)}_{\text{minimized by } q = \pi} \right\}, \quad (2)$$

Notation:

- $\mathcal{P}(\Theta)$ = all probability distributions on Θ
- 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 θ for x_i
- Inference = regularizing MLE $\hat{\theta}_n$ with $\text{KLD}(q||\pi)$

(2.2) Modularity and Interpretation

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

3.) Reduce computational cost by optimizing over a smaller set $\Pi \subset \mathcal{P}(\Theta)$
(= Variational Inference!)

2.) Guard against likelihood misspecification by using robust model scoring rules $\mathcal{L}(\cdot)$ instead of $-\log(\cdot)$

1.) Suppress ill-informed priors by choosing a weaker regularizing divergence $D(\cdot||\pi)$ instead of $\text{KLD}(\cdot||\pi)$

(2.3) A natural generalization

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

An optimization-centric generalization of Bayesian inference:

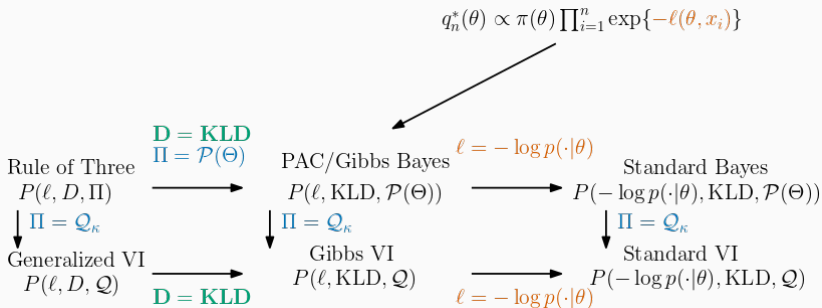
- (1) $D(\cdot \parallel \pi)$ = **any divergence regularizer** penalizing deviations from π
- (2) $\ell(\theta, \mathbf{x})$ = **any loss** assessing how well θ and \mathbf{x} fit together
- (3) $\Pi \subseteq \mathcal{P}(\Theta)$ = **some subset** of all probability distributions on Θ

\implies **Shorthand Notation:** $P(\ell, D, \Pi)$

(2.4) Special cases of note

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

(Some) special cases of interest:



(2.5) Generalized Variational Inference

Terminology:

- Presented so far: a conceptual generalization of Bayesian inference
- **Generalized Variational Inference (GVI)** = computable inference algorithms based on that generalization
- \implies means that $\Pi = \mathcal{Q}$ for some variational family \mathcal{Q} .

Open questions:

- General properties of this generalisation?
- Applications/use cases?

\implies answered next

(3) General features of the generalization

Answers to three questions:

- (3.1) 'Is this even reasonable?' (sanity checks)
- (3.2) 'Okay, but under which conditions is it reasonable?' (axiomatic justification)
- (3.3) 'Fine, but can it help me understand existing methods?' (VI, PAC-Bayes bounds, power posteriors, ...)

(3.1) 'Sanity checks': existence, uniqueness, consistency

$$q_n^*(\theta) = \arg \min_{q \in \Pi} \left\{ \mathbb{E}_{q(\theta)} \left[\sum_{i=1}^n \ell(\theta, x_i) \right] + \mathbf{D}(q||\pi) \right\} = P(\ell, \mathbf{D}, \Pi)$$

Existence: If Π is convex and chosen so that $q \mapsto \mathbb{E}_{q(\theta)} [\sum_{i=1}^n \ell(\theta, x_i)]$ and $q \mapsto \mathbf{D}(q||\pi)$ are continuous on Π , then the minimum exists whenever $q \mapsto \mathbf{D}(q||\pi)$ is convex.

\implies basic convex analysis on Banach spaces

Uniqueness: Guaranteed if q_n^* exists and $q \mapsto \mathbf{D}(q||\pi)$ is *strictly* convex

\implies basic convex analysis on Banach spaces

(3.1) 'Sanity checks': existence, uniqueness, consistency

$$q_n^*(\theta) = \arg \min_{q \in \Pi} \left\{ \mathbb{E}_{q(\theta)} \left[\sum_{i=1}^n \ell(\theta, x_i) \right] + D(q \parallel \pi) \right\} = P(\ell, D, \Pi)$$

Beyond continuity & convexity of $q \mapsto D(q \parallel \pi)$? You can still show existence and uniqueness without convex regularizers [need that losses are norm-coercive or Θ is compact; Arguments not that basic; see Lemma 1 in Knoblauch (2019)]

Consistency: Guaranteed under relatively mild regularity conditions; arguments are unfortunately quite complicated and rely on Γ -convergence, see Knoblauch (2019).

\implies Idea of Γ -convergence \approx 'if a sequence of objectives Γ -converge, then their minimizers also converge (in a suitable sense, usually weakly)'

(3.2) Axiomatic justification

Q: Does this generalisation result from 'reasonable' axioms?

Axiom 1 (Variational representation)

The posterior $q^* \in \mathcal{P}(\Theta)$ solves an optimization problem over some space $\Pi \subseteq \mathcal{P}(\Theta)$. For any finite sample $\{x_i\}_{i=1}^n$, the optimization problem seeks to jointly minimize two criteria:

- (i) An 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 some statistical divergence D .

Theorem 1 (Form 1)

Under Axiom 1, posterior belief distributions can be written as

$$q^*(\theta) = \arg \min_{q \in \Pi} \left\{ f \left(\mathbb{E}_{q(\theta)} \left[\sum_{i=1}^n \ell(\theta, x_i) \right], D(q \parallel \pi) \right) \right\},$$

where $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is some function that may depend on $\pi, \Pi, \ell, \{x_i\}_{i=1}^n$, or D .

(3.2) Axiomatic justification

Axiom 2 (Recovers Bayesian Posteriors)

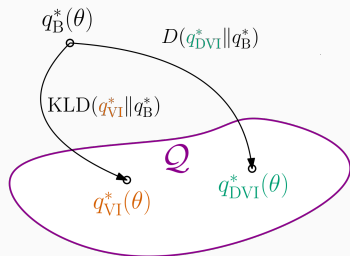
Function f in Theorem 1 does not depend on $\pi, \Pi, \ell, \{x_i\}_{i=1}^n$, or D .
Further, q^* is the Gibbs posterior if $D = \text{KLD}$, $\Pi = \mathcal{P}(\Theta)$.

Theorem 2

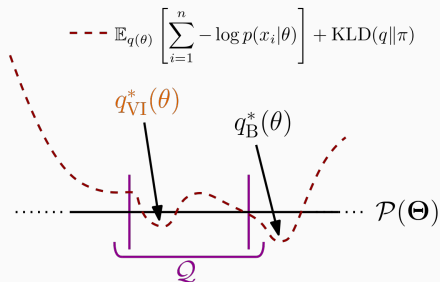
Suppose the posterior belief $q^ \in \mathcal{P}(\Theta)$ satisfies Axioms 1 and 2. Then the objective of Theorem 1 is uniquely identified as $f(x, y) = x + y$ so that*

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

(3.3) Relation to VI / PAC-Bayes / power posteriors



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



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

(3.3) Relation to VI / PAC-Bayes / power posteriors

Proposition 1 (Optimality of standard VI)

Relative to the infinite-dimensional optimization problem over $\mathcal{P}(\Theta)$ characterizing the Gibbs posterior and a fixed variational family Π , standard VI produces the optimal solution (i.e. posterior belief) in Π .

Proof.

VI posteriors are minima of the same objective as the full Bayesian posterior—but constrained to some subset Π . □

Proposition 2 (Suboptimality of alternative methods)

Relative to the infinite-dimensional problem over $\mathcal{P}(\Theta)$ characterizing Gibbs posteriors, and relative to a fixed finite-dimensional variational family Π , non-standard VI methods produce sub-optimal solutions (i.e., posterior beliefs).

(3.3) Relation to VI / PAC-Bayes / power posteriors

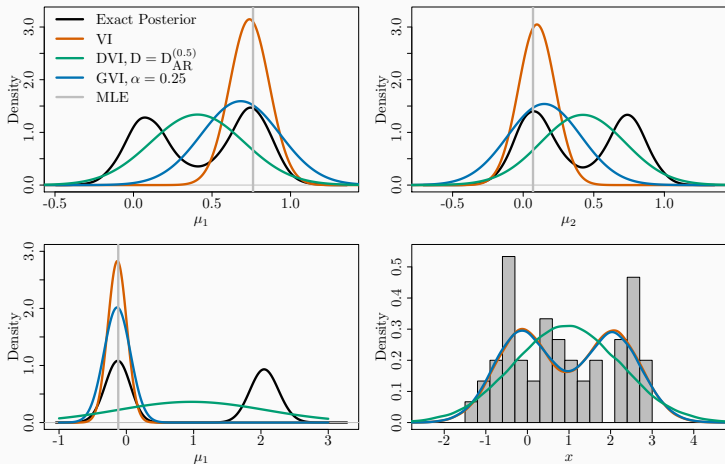


Figure 2 – Top row: marginal posteriors for location parameters μ_1, μ_2 in 2D Bayesian Gaussian Mixture Model. **Bottom left:** same marginal for μ_1 as we increase $|\mu_1 - \mu_2|$. **Bottom right:** posterior predictives.

(3.3) Relation to VI / PAC-Bayes / power posteriors

E.g., McAllester's (original) PAC-Bayes bound: if $x_i \stackrel{iid}{\sim} \mu$ so that true risk is $R(\theta) = \mathbb{E}_{x \sim \mu}[\ell(\theta, x)]$ and $a \leq \ell \leq b$, then uniformly for all $q \in \mathcal{P}(\Theta)$ with probability at least $1 - \varepsilon$,

$$\mathbb{E}_{q(\theta)} [R(\theta)] \leq \mathbb{E}_{q(\theta)} \left[\frac{1}{n} \sum_{i=1}^n \ell(\theta, x_i) \right] + \sqrt{\frac{\text{KLD}(q, \pi) + \log \frac{2\sqrt{n}}{\varepsilon}}{2n}}.$$

I.e., we can minimize the righthand side over $q \in \mathcal{P}(\Theta)$ to find that **the tightest generalisation bound** has a solution of the form

$P(\ell, D_{MCA}, \mathcal{P}(\Theta))$ with

$$D_{MCA}(q \parallel \pi) = \sqrt{n} \cdot \left(\sqrt{\frac{\text{KLD}(q, \pi) + \log \frac{2\sqrt{n}}{\varepsilon}}{2}} - \sqrt{\frac{\log \frac{2\sqrt{n}}{\varepsilon}}{2}} \right).$$

\implies PAC-Bayes [!] a way to justify non-standard prior regularization, e.g. Bégin et al. (2016) [Rényi's α -divergence], Alquier & Guedj (2018), Ohnishi & Honorio (2020) [f -divergences]

(3.3) Relation to VI / PAC-Bayes / power posteriors

Power posteriors/'cold posteriors': For some $\beta > 0$, given by

$$\begin{aligned} q^*(\theta) &\propto \pi(\theta) \prod_{i=1}^n p(x_i|\theta)^\beta = P(-\log p(\cdot|\theta) \cdot \beta, \text{KLD}, \mathcal{P}(\Theta)) \quad (3) \\ &= P(-\log p(\cdot|\theta), \text{KLD} \cdot \frac{1}{\beta}, \mathcal{P}(\Theta)) \end{aligned}$$

Cold posteriors: $\beta > 1$, i.e. more weight on data rather than prior
 \implies often used for NNs, where our priors are extremely poor.

Power posteriors: $\beta < 1$, i.e. more weight on prior rather than data
 \implies often used when likelihoods are poorly misspecified and advertised as 'robust' (... it's always more robust as $n \rightarrow \infty$, but if n is small, your robustness properties will depend a lot on the prior)

(3.3) Relation to VI / PAC-Bayes / power posteriors / ...

Method	$\ell(\boldsymbol{\theta}, x_j)$	D	Π
Standard Bayes	$-\log p(x_j \boldsymbol{\theta})$	KLD	$\mathcal{P}(\Theta)$
Power Likelihood Bayes	$-\log p(x_j \boldsymbol{\theta})$	$\frac{1}{w}$ KLD, $w < 1$	$\mathcal{P}(\Theta)$
Composite Likelihood Bayes	$-w_j \log p(x_j \boldsymbol{\theta})$	KLD	$\mathcal{P}(\Theta)$
Divergence-based Bayes	divergence-based ℓ	KLD	$\mathcal{P}(\Theta)$
PAC/Gibbs Bayes	any ℓ	any D	$\mathcal{P}(\Theta)$
VAE	$-\log p_{\zeta}(x_j \boldsymbol{\theta})$	KLD	\mathcal{Q}
β -VAE	$-\log p_{\zeta}(x_j \boldsymbol{\theta})$	$\beta \cdot$ KLD, $\beta > 1$	\mathcal{Q}
Bernoulli-VAE	continuous Bernoulli	KLD	\mathcal{Q}
Standard VI	$-\log p(x_j \boldsymbol{\theta})$	KLD	\mathcal{Q}
Power VI	$-\log p(x_j \boldsymbol{\theta})$	$\frac{1}{w}$ KLD, $w < 1$	\mathcal{Q}
Utility VI	$-\log p(x_j \boldsymbol{\theta}) + \log u(h, x_j)$	KLD	\mathcal{Q}
Regularized Bayes	$-\log p(x_j \boldsymbol{\theta}) + \phi(\boldsymbol{\theta}, x_j)$	KLD	\mathcal{Q}
Gibbs VI	any ℓ	KLD	\mathcal{Q}
Generalized VI	any ℓ	any D	\mathcal{Q}

(3.4) Closed Forms

Well-known: Gibbs posteriors

$$P(\ell, \text{KLD}, \mathcal{P}(\Theta)) \propto \pi(\theta) \cdot \exp \left\{ - \sum_{i=1}^n \ell(\theta, x_i) \right\}$$

Unknown (until a few months ago): What if $D \neq \text{KLD}$?

\implies We now know the general form if D an ϕ -divergence

Proposition 3.1. Assume that ϕ is differentiable, strictly convex and define $\tilde{\phi}$ on \mathbb{R} by $\tilde{\phi}(x) = \phi(x)$ if $x \geq 0$ and $\tilde{\phi}(x) = +\infty$ otherwise. Then

$$\tilde{\phi}^* = \sup_{x \in \mathbb{R}} [xy - \tilde{\phi}(x)] = \sup_{x \geq 0} [xy - \phi(x)] \quad (3.1)$$

is differentiable and for any $y \in \mathbb{R}$,

$$\nabla \tilde{\phi}^*(y) = \operatorname{argmax}_{x \geq 0} \{xy - \phi(x)\}. \quad (3.2)$$

Assume moreover than $\tilde{\phi}^*(\lambda - a) - \lambda \rightarrow \infty$ when $\lambda \rightarrow \infty$, for any $a \geq 0$. Then

$$\lambda_t \in \operatorname{argmin}_{\lambda \in \mathbb{R}} \left\{ \int \tilde{\phi}^* \left(\lambda - \eta \sum_{s=1}^{t-1} \ell_s(\theta) \right) \pi(d\theta) - \lambda \right\} \quad (3.3)$$

exists, and

$$\rho^t(d\theta) = \nabla \tilde{\phi}^* \left(\lambda_t - \eta \sum_{s=1}^{t-1} \ell_s(\theta) \right) \pi(d\theta) \quad (3.4)$$

minimizes (1.4).

(3.4) Closed Forms

Two examples:

Example 3.2 (χ^2 -divergence). We come back to the example $\phi(x) = x^2 - 1$, $D_\phi(\rho||\pi) = \chi^2(\rho||\pi)$ the chi-squared divergence. In this case, $\phi^*(y) = (y^2/4)\mathbf{1}_{\{y \geq 0\}}$ and so $\nabla \tilde{\phi}^*(y) = (y/2)_+$. This leads to

$$\rho^t(d\theta) = \left[\frac{\lambda_t - \eta \sum_{s=1}^{t-1} \ell_s(\theta)}{2} \right]_+ \pi(d\theta). \quad (3.7)$$

In this case, λ_t is not available in closed form, but it is the only constant that will make the above sum to 1.

Example 3.3 (p -power divergence). More generally, consider $\phi(x) = x^p - 1$. In this case $\nabla \tilde{\phi}^*(y) = (y/p)_+^{1/(p-1)}$. This leads to

$$\rho^t(d\theta) = \left[\frac{\lambda_t - \eta \sum_{s=1}^{t-1} \ell_s(\theta)}{p} \right]_+^{\frac{1}{p-1}} \pi(d\theta). \quad (3.8)$$

(4) Applications

Problems that can be tackled:

- Robustness to poor priors
- Robustness to poor likelihoods
- Simplified computation
- ...

(4.1) Robustness to prior misspecification

$$q^*(\theta) = \arg \min_{q \in \Pi} \left\{ \mathbb{E}_{q(\theta)} \left[\sum_{i=1}^n \ell(\theta, x_i) \right] + D(q||\pi) \right\} = P(\ell, \mathbf{D}, \Pi)$$

What we want: D that behaves like **KLD** if π is reasonable, but ignores it if the data don't fit the prior at all

First idea: down-weight $D = \text{KLD}$ like in cold posteriors

\implies Problem: Now, not being certain of your prior amounts to being more certain in your posterior....

Q: Is there an alternative?

\implies For reasons we don't understand fully[†], Rényi's α -divergence seems to do behave exactly as we want! (small loss of efficiency observed if prior is well-specified)

[†]some limited theoretical results in Theorem 14 of Knoblauch et al. (2019)

(4.1) Robustness to prior misspecification

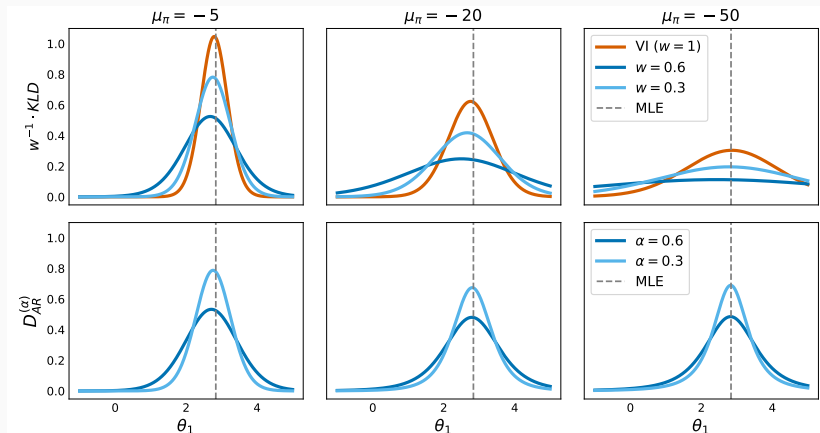


Figure 3 – The prior for the coefficients is a Normal Inverse Gamma distribution given by $\mu \sim \mathcal{NI}^{-1}(\mu_{\pi} \cdot \mathbf{1}_d, v_{\pi} \cdot \mathbf{I}_d, a_{\pi}, b_{\pi})$ with $v_{\pi} = 4 \cdot \mathbf{I}_d$, $a_{\pi} = 3$, $b_{\pi} = 5$ and various values for μ_{π} . For all posteriors, the loss ℓ is the correctly specified negative log likelihood; all posteriors lie inside a mean field normal family \mathcal{Q} .

(4.1) Robustness to prior misspecification

Example: Bayesian Neural Networks

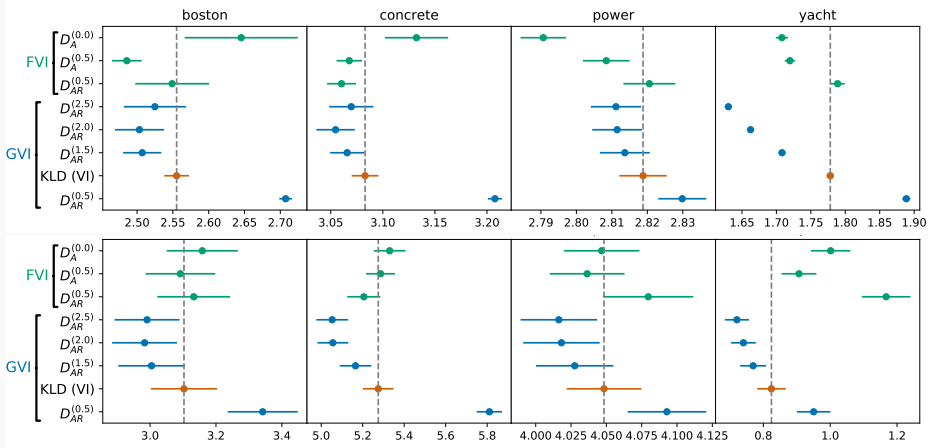


Figure 4 – Top: Negative test log likelihoods. Bottom row: Test RMSE.

(4.2) Robustness to likelihood misspecification

What it means: [Taking definition from Hooker & Vidyashankar]

Consider ε -contamination model of size $\varepsilon \in (0, 1)$

$$\mathbb{P}_{n,\varepsilon,y} = (1 - \varepsilon)\mathbb{P}_n + \varepsilon\delta_y; \quad y \in \mathcal{X}$$

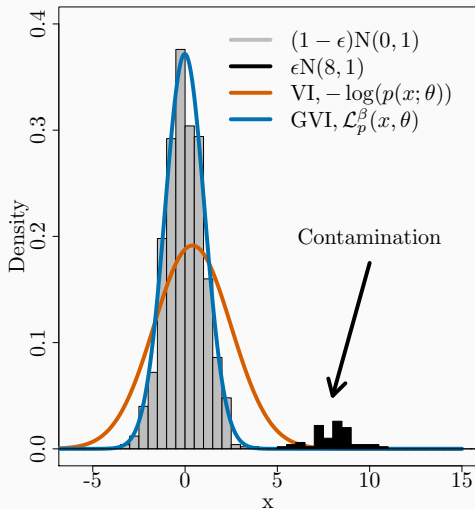
and write $\ell_{\varepsilon,n}(\theta) = \mathbb{E}_{x \sim \mathbb{P}_{n,\varepsilon,y}}[\ell(\theta, x)]$. Define the *posterior influence function* as

$$\text{PIF}(y, \theta, \mathbb{P}_n) := \frac{d}{d\varepsilon} P(\ell_{\varepsilon,n}(\theta), \mathbf{D}, \mathbf{\Pi})|_{\varepsilon=0}.$$

The posterior $P(\ell_{\varepsilon,n}(\theta), \mathbf{D}, \mathbf{\Pi})$ is called *globally bias-robust* if $\sup_{\theta \in \Theta} \sup_{y \in \mathcal{X}} |\text{PIF}(y, \theta, \mathbb{P}_n)| < \infty$, meaning that the sensitivity of the generalised posterior to the contaminant y is limited.

(4.2) Robustness to likelihood misspecification

What it means:



(4.2) Motivation for Discrepancy-based losses

General Observation: log likelihoods \approx minimize the KLD

$$\begin{aligned}q_n^*(\theta) &= \arg \min_{q \in \Pi} \left\{ \mathbb{E}_{q(\theta)} \left[\sum_{i=1}^n -\log p(x_i|\theta) \right] + \text{KLD}(q||\pi) \right\} \\&= \arg \min_{q \in \Pi} \left\{ \mathbb{E}_{q(\theta)} \left[\frac{1}{n} \sum_{i=1}^n -\log \frac{p(x_i|\theta)}{p_0(x_i)} \right] - \frac{1}{n} \sum_{i=1}^n \log p_0(x_i) + \frac{1}{n} \text{KLD}(q||\pi) \right\} \\&= \arg \min_{q \in \Pi} \left\{ \underbrace{\mathbb{E}_{q(\theta)} \left[\frac{1}{n} \sum_{i=1}^n -\log \frac{p(x_i|\theta)}{p_0(x_i)} \right]}_{\approx \text{KLD}(p_0||p(\cdot|\theta))} + \frac{1}{n} \text{KLD}(q||\pi) \right\}\end{aligned}$$

Obvious question: What are the (dis)advantages of minimizing other discrepancies between p_0 and $p(\cdot|\theta)$ instead?

General Answer: Decreases statistical efficiency, increases robustness

(4.2) Robustness to likelihood misspecification

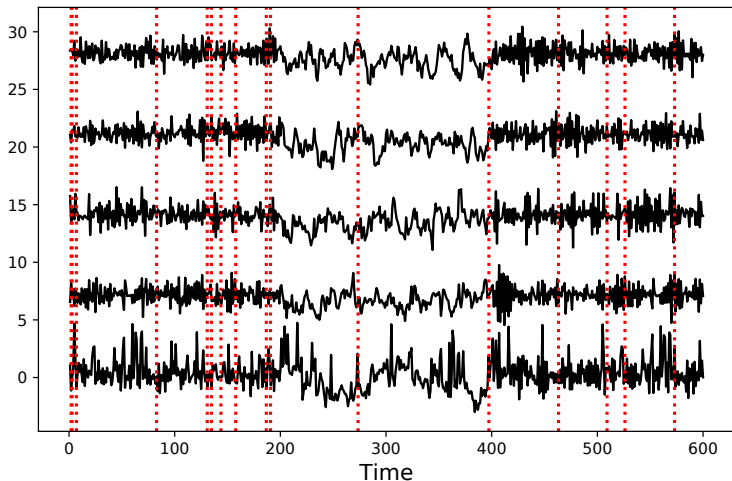
Some applications of robust divergences as losses:

- Reducing false detection of changepoints
- Improving performance of deep Gaussian Processes
- Graphical models
- ...

Many, many more... — worth considering every time your likelihood is at best a reasonable guess.

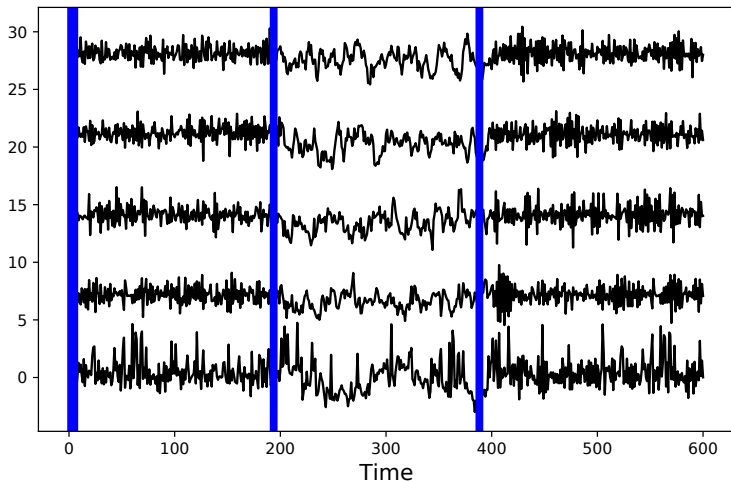
(4.2) Outlier-Robust Changepoint Detection

Using standard Bayesian On-line Changepoint Detection



(4.2) Outlier-Robust Changepoint Detection

Using the β -divergence for Robust Changepoint Detection



(4.2) Likelihoods with Deep Gaussian Processes

γ -divergences for DGP regression³

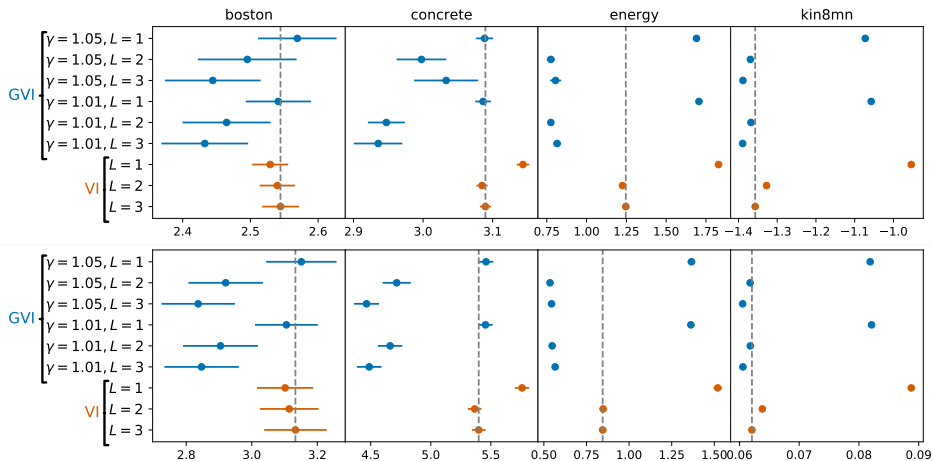


Figure 5 – Top row: Negative test log likelihoods. Bottom row: Test RMSE.

(4.3) Improving Bayesian computation

Idea: Find losses $\tilde{\ell} \neq \ell$ s.t.

$$\text{CompTime} [P(\tilde{\ell}, \mathbf{D}, \Pi)] \ll \text{CompTime} [P(\ell, \mathbf{D}, \Pi)]$$

Example 1: Approximate Bayesian Computation (ABC).

Step 1: $\theta \sim \pi(\theta)$

Step 2: sample $x_{\text{fake}} \sim p(x_{\text{fake}}|\theta)$

Step 3: keep θ if $D(x_{\text{fake}}, x_{\text{observed}}) < \varepsilon$; discard otherwise.

\implies

$$q(\theta|x_{\text{observed}}) \approx q_{\text{abc}}(\theta|x_{\text{observed}}) \propto \int \mathbf{1}_{[D(x_{\text{fake}}, x_{\text{observed}}) < \varepsilon]} p(x_{\text{fake}}|\theta) \pi(\theta) dx_{\text{fake}}.$$

Usually: $q_{\text{abc}}(\theta|x_{\text{observed}})$ interpreted as approximation to $q(\theta|x_{\text{observed}})$.

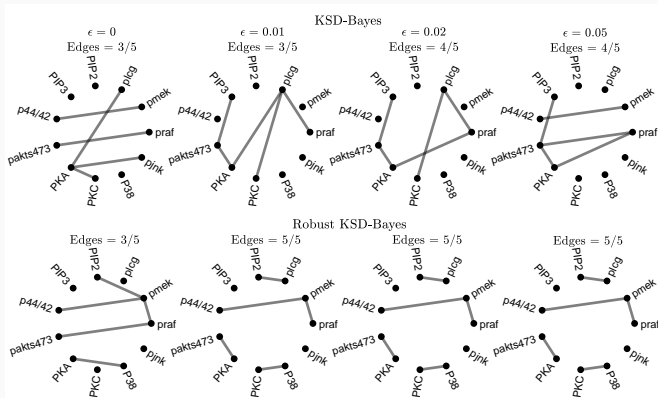
Alternative: $q_{\text{abc}}(\theta|x_{\text{observed}}) \propto \pi(\theta) \exp\left\{-\underbrace{L_n(\theta, x_{\text{observed}})}\right\}$

$$= \int \mathbf{1}_{[D(x_{\text{fake}}, x_{\text{observed}}) < \varepsilon]} p(x_{\text{fake}}|\theta)$$

\implies simulator error model $err(x_{\text{fake}}, x_{\text{observed}}) = \int \mathbf{1}_{[D(x_{\text{fake}}, x_{\text{observed}}) < \varepsilon]}$

\implies 'smoother' error models more sample-efficient (e.g., $err = \text{normal}$)

(4.3) Improving Bayesian computation



(4.3) Improving Bayesian computation

Idea: Find losses $\tilde{\ell} \neq \ell$ s.t.

CompTime $[P(\tilde{\ell}, \mathbf{D}, \Pi)] \ll \text{CompTime} [P(\ell, \mathbf{D}, \Pi)]$

Example 2: Intractable Likelihoods/Energy-based models

Challenge: Likelihoods with unknown normalisers, i.e.

$$p(x|\theta) = \underbrace{\hat{p}(x|\theta)}_{\text{known}} \cdot \underbrace{Z(\theta)}_{\text{unknown}}$$

\implies standard Bayesian posterior 'doubly intractable'

Solution: loss depending on $p(x|\theta)$ *only* via $\nabla_x p(x|\theta) \stackrel{!}{=} \nabla_x \hat{p}(x|\theta)$

\implies Stein's method! Operationalisable via Kernel Stein Discrepancies

\implies Makes posteriors 'singly intractable'

... **Question:** but can we do even better?!

(4.3) Improving Bayesian computation

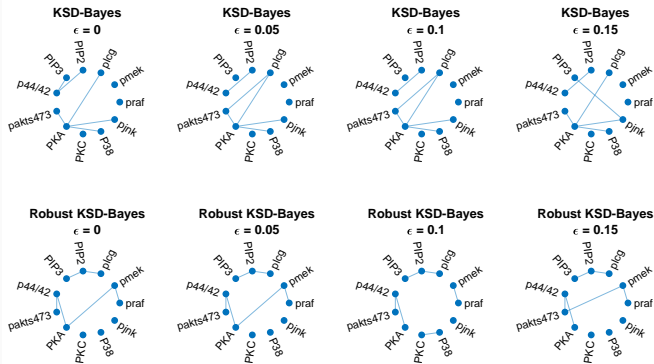
Answer: Yes, whenever $p(x|\theta)$ is part of the exponential family!

⇒ Then, we get **closed forms** if π is normal!

⇒ Closed forms instead of 'doubly intractable' posteriors

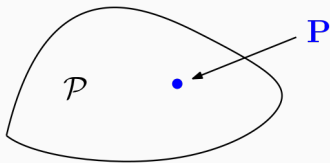
⇒ Added bonus: robustness to model misspecification!

$$p(x|\theta) \propto \exp\left(-\sum_i \theta_i \exp(x_i) - \sum_{i < j} \theta_{i,j} \exp(x_i) \exp(x_j)\right) \times \exp\left(\sum_i x_i\right)$$



Summary

- (1) Bayesian inference = a **single** optimization problem $P \in \mathcal{P}$
- (2) Optimization problem $P \in \mathcal{P}$ needs **strict assumptions**
- (3) These assumptions are often violated in Machine Learning (ML)
- (4) **Many** other optimization problems $P' \in \mathcal{P}$ you could solve instead (Even though we still know relatively little about them!)
- (5) In ML, we would often do better to solve $P' \in \mathcal{P}$ rather than P



Also—And perhaps even more importantly:

- The study of $P' \in \mathcal{P}$ has just begun! Get involved! :)
- If I managed to inspire you to work on these problems, get in touch! I love to collaborate, and there are enough open problems for a lifetime! :)

Contact me/follow my research at

E-mail address: j.knoblach@warwick.ac.uk

Personal Website: <https://jeremiasknoblach.github.io/>

Twitter Handle: [@LauchLab](https://twitter.com/LauchLab)

