A Regret Bound for Online Variational Inference

Pierre Alquier **Pierre Alquier Pierre Alquie**

Heilbronn Institute - University of Bristol - Dec. 5, 2019

Pierre Alquier, RIKEN AIP Online variational inference



Badr-Eddine Chérief-Abdellatif



ENSAE

IP PARIS

Emtiyaz Khan







https://emtiyaz.github.io/

Pierre Alquier, RIKEN AIP Online variational inference



K. Osawa, S. Swaroop, A. Jain, R. Eschenhagen, R. E. Turner, R. Yokota, M. E. Khan (2019). *Practical Deep Learning with Bayesian Principles*. NeurIPS.

K. Osawa, S. Swaroop, A. Jain, R. Eschenhagen, R. E. Turner, R. Yokota, M. E. Khan (2019). *Practical Deep Learning with Bayesian Principles.* NeurIPS.

- proposes a fast algorithm to approximate the posterior,
- applies it to train Deep Neural Networks on CIFAR-10, ImageNet ...
- observation : improved uncertainty quantification.



K. Osawa, S. Swaroop, A. Jain, R. Eschenhagen, R. E. Turner, R. Yokota, M. E. Khan (2019). Practical Deep Learning with Bayesian Principles. NeurIPS.

- proposes a fast algorithm to approximate the posterior,
- applies it to train Deep Neural Networks on CIFAR-10, ImageNet ...
- observation : improved uncertainty quantification.



Objective : provide a theoretical analysis of this algorithm.

K. Osawa, S. Swaroop, A. Jain, R. Eschenhagen, R. E. Turner, R. Yokota, M. E. Khan (2019). *Practical Deep Learning with Bayesian Principles*. NeurIPS.

- proposes a fast algorithm to approximate the posterior,
- applies it to train Deep Neural Networks on CIFAR-10, ImageNet ...
- observation : improved uncertainty quantification.



Objective : provide a theoretical analysis of this algorithm. **First step** : simplified versions.

Sequential prediction problem

Sequential prediction problem

1 1 x_1 given

Sequential prediction problem



2 predict $y_1 : \hat{y}_1$

Sequential prediction problem

- **1 1** x_1 given
 - **2** predict $y_1 : \hat{y}_1$
 - **3** y_1 is revealed

Sequential prediction problem

x₁ given
predict y₁ : ŷ₁
y₁ is revealed
x₂ given

Sequential prediction problem

x₁ given
predict y₁ : ŷ₁
y₁ is revealed
x₂ given
predict y₂ : ŷ₂

Sequential prediction problem

x₁ given
predict y₁ : ŷ₁
y₁ is revealed
x₂ given
predict y₂ : ŷ₂
y₂ revealed

Sequential prediction problem

x₁ given
predict y₁ : ŷ₁
y₁ is revealed
x₂ given
predict y₂ : ŷ₂
y₂ revealed
x₃ given

Sequential prediction problem



Sequential prediction problem



Sequential prediction problem			
0	0 2 3	x_1 given predict y_1 : \hat{y}_1 y_1 is revealed	Objective :
2	0 0	x_2 given predict y_2 : \hat{y}_2 y_2 revealed	
3	0 0	x_3 given predict y_3 : \hat{y}_3 y_3 revealed	
❹.			

Sequential prediction problem



Objective : make sure that we learn to predict well **as soon as possible**.



Objective : make sure that we learn to predict well **as soon as possible**. Keep



as small as possible, without unrealistic assumptions on the data.

As an example, consider for one minute a classification problem, that is :

$$y_t \in \{0,1\}$$
 and $\ell(\hat{y}_t, y_t) = \mathbf{1}_{\{y_t \neq \hat{y}_t\}}.$

As an example, consider for one minute a classification problem, that is :

$$y_t \in \{0,1\} \text{ and } \ell(\hat{y}_t, y_t) = \mathbf{1}_{\{y_t \neq \hat{y}_t\}}.$$

In the worst case scenario, y_t generated by an omniscient opponent by : $y_t = 1 - \hat{y}_t$. Then

$$\sum_{t=1}^T \ell(\hat{y}_t, y_t) = T.$$

As an example, consider for one minute a classification problem, that is :

$$y_t \in \{0,1\} \text{ and } \ell(\hat{y}_t, y_t) = \mathbf{1}_{\{y_t \neq \hat{y}_t\}}.$$

In the worst case scenario, y_t generated by an omniscient opponent by : $y_t = 1 - \hat{y}_t$. Then

$$\sum_{t=1}^T \ell(\hat{y}_t, y_t) = T.$$

On the other hand, many real world phenomena can be "quite well" described by models. These models allow to do "sensible" predictions.

As an example, consider for one minute a classification problem, that is :

$$y_t \in \{0,1\} \text{ and } \ell(\hat{y}_t, y_t) = \mathbf{1}_{\{y_t \neq \hat{y}_t\}}.$$

In the worst case scenario, y_t generated by an omniscient opponent by : $y_t = 1 - \hat{y}_t$. Then

$$\sum_{t=1}^T \ell(\hat{y}_t, y_t) = T.$$

On the other hand, many real world phenomena can be "quite well" described by models. These models allow to do "sensible" predictions.

The extreme case would be the constraint $y_t = f(x_t)$, where $f \in \mathcal{F}$ for a known class \mathcal{F} – the *realizable case*. Let's study it as a toy example when \mathcal{F} is finite.

A naive strategy

Here
$$y_t = f_{i^*}(x_t)$$
 where $i^* \in \{1, \ldots, M\}$ is unknown.

Naive strategy

Here
$$y_t = f_{i^*}(x_t)$$
 where $i^* \in \{1, \ldots, M\}$ is unknown.

Here
$$y_t = f_{i^*}(x_t)$$
 where $i^* \in \{1, \ldots, M\}$ is unknown.

• predict
$$\hat{y}_t = f_{i(t)}(x_t)$$
, observe y_t ,

Here
$$y_t = f_{i^*}(x_t)$$
 where $i^* \in \{1, \ldots, M\}$ is unknown.

Here
$$y_t = f_{i^*}(x_t)$$
 where $i^* \in \{1, \ldots, M\}$ is unknown.

Start with i(1) = 1 and $C(1) = \{1, ..., M\}$. At step t,

Theorem

$$\forall T, \sum_{t=1}^{T} \ell(\hat{y}_t, y_t) \leq M-1.$$

(Still
$$y_t = f_{i^*}(x_t)$$
 where $i^* \in \{1, \dots, M\}$ is unknown).

The halving algorithm

(Still
$$y_t = f_{i^*}(x_t)$$
 where $i^* \in \{1, \ldots, M\}$ is unknown).

The halving algorithm

(Still
$$y_t = f_{i^*}(x_t)$$
 where $i^* \in \{1, \ldots, M\}$ is unknown).

The halving algorithm

Start with i(1) = 1 and $C(1) = \{1, ..., M\}$. At step *t*,

• predict $\hat{y}_t =$ "majority vote in C(t)", observe y_t ,

(Still
$$y_t = f_{i^*}(x_t)$$
 where $i^* \in \{1, \ldots, M\}$ is unknown).

The halving algorithm

Start with i(1) = 1 and $C(1) = \{1, ..., M\}$. At step *t*,

• predict $\hat{y}_t =$ "majority vote in C(t)", observe y_t ,

2 update
$$C(t+1) = \{i \in C(t) : f_i(x_t) = y_t\}.$$

(Still
$$y_t = f_{i*}(x_t)$$
 where $i^* \in \{1, \ldots, M\}$ is unknown).

The halving algorithm

Start with i(1) = 1 and $C(1) = \{1, ..., M\}$. At step *t*,

• predict $\hat{y}_t =$ "majority vote in C(t)", observe y_t ,

2 update
$$C(t+1) = \{i \in C(t) : f_i(x_t) = y_t\}.$$

Theorem

$$orall \mathcal{T}$$
, $\sum_{t=1}^{T} \ell(\hat{y}_t, y_t) \leq \log_2(M)$.

Two extremes :

- playing against the devil $y_t = 1 \hat{y}_t$,
- assuming a true, exact model \mathcal{F} .

Two extremes :

- playing against the devil $y_t = 1 \hat{y}_t$,
- assuming a true, exact model \mathcal{F} .

Real-life is somewhere in between !

Two extremes :

- playing against the devil $y_t = 1 \hat{y}_t$,
- assuming a true, exact model \mathcal{F} .

Real-life is somewhere in between !

Objective

Strategy such that

$$\sum_{t=1}^{T} \ell(\hat{y}_t, y_t) \leq \inf_{\substack{f \in \mathcal{F} \\ t=1 \\ = \text{ T in the worst case (devil),} \\ = 0 in the ideal case (true model),}}^{T} \ell(f(x_t), y_t) + \underbrace{\mathcal{B}(T)}_{\text{as small as possible !!}}.$$
Online gradient algorithm (OGA)

Given

- a set of predictors $\{f_{\theta}, \theta \in \Theta \subset \mathbb{R}^d\}$, e.g $f_{\theta}(x) = \langle \theta, x \rangle$,
- an initial guess θ_1 ,

$$\hat{y}_t = f_{\theta_t}(x_t) \text{ and } \theta_{t+1} = \theta_t - \eta \nabla_{\theta} \ell(f_{\theta_t}(x_t), y_t).$$

Online gradient algorithm (OGA)

Given

- a set of predictors $\{f_{\theta}, \theta \in \Theta \subset \mathbb{R}^d\}$, e.g $f_{\theta}(x) = \langle \theta, x \rangle$,
- an initial guess θ_1 ,

$$\hat{y}_t = f_{\theta_t}(x_t)$$
 and $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \ell_t(\theta_t)$.

Online gradient algorithm (OGA)

Given

- a set of predictors $\{f_{\theta}, \theta \in \Theta \subset \mathbb{R}^d\}$, e.g $f_{\theta}(x) = \langle \theta, x \rangle$,
- an initial guess θ_1 ,

$$\hat{y}_t = f_{\theta_t}(x_t)$$
 and $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \ell_t(\theta_t)$.

Note that θ_{t+1} can be obtained by :

$$\min_{\theta} \left\{ \left\langle \theta, \nabla_{\theta} \ell_t(\theta_t) \right\rangle + \frac{\|\theta - \theta_t\|^2}{2\eta}, \right\}$$

$$\min_{\theta} \left\{ \left\langle \theta, \sum_{s=1}^t \nabla_{\theta} \ell_s(\theta_s) \right\rangle + \frac{\|\theta - \theta_1\|^2}{2\eta} \right\}.$$

Guarantees for OGA

Theorem - classical, see Shalev-Shwarz 2012

Assume ℓ is convex and L-Lipschitz. Then

$$\sum_{t=1}^{T} \ell_t(\theta_t) \leq \inf_{\theta \in \Theta} \left\{ \sum_{t=1}^{T} \ell_t(\theta) + \eta T L^2 + \frac{\|\theta\|^2}{2\eta} \right\}.$$

S. Shalev-Shwartz (2012). Online learning and online convex optimization. Foundations and Trends in Machine Learning.

Guarantees for OGA

Theorem - classical, see Shalev-Shwarz 2012

Assume ℓ is convex and *L*-Lipschitz. Then

$$\sum_{t=1}^{T} \ell_t(\theta_t) \leq \inf_{\theta \in \Theta} \left\{ \sum_{t=1}^{T} \ell_t(\theta) + \eta T L^2 + \frac{\|\theta\|^2}{2\eta} \right\}.$$

S. Shalev-Shwartz (2012). Online learning and online convex optimization. Foundations and Trends in Machine Learning.

$$\sum_{t=1}^{T} \ell_t(\theta_t) \leq \inf_{\|\theta\| \leq B} \sum_{t=1}^{T} \ell_t(\theta) + \eta T L^2 + \frac{B^2}{2\eta}$$

Guarantees for OGA

Theorem - classical, see Shalev-Shwarz 2012

Assume ℓ is convex and *L*-Lipschitz. Then

$$\sum_{t=1}^{T} \ell_t(\theta_t) \leq \inf_{\theta \in \Theta} \left\{ \sum_{t=1}^{T} \ell_t(\theta) + \eta T L^2 + \frac{\|\theta\|^2}{2\eta} \right\}.$$

S. Shalev-Shwartz (2012). Online learning and online convex optimization. Foundations and Trends in Machine Learning.

$$\sum_{t=1}^{T} \ell_t(\theta_t) \leq \inf_{\|\theta\| \leq B} \sum_{t=1}^{T} \ell_t(\theta) + \eta T L^2 + \frac{B^2}{2\eta}$$

The choice $\eta = B/(L\sqrt{2T})$ leads to

$$\sum_{t=1}^{T} \ell_t(\theta_t) \leq \inf_{\|\theta\| \leq B} \sum_{t=1}^{T} \ell_t(\theta) + BL\sqrt{2T}.$$

Bayesian learning and variational inference (VI)

$$\pi_{t+1}(\theta) := \pi(\theta|x_1, y_1, \dots, x_t, y_t) \propto \exp\left(-\eta \sum_{s=1}^t \ell_s(\theta)\right) \pi(\theta).$$

Bayesian learning and variational inference (VI)

$$\pi_{t+1}(\theta) := \pi(\theta|x_1, y_1, \dots, x_t, y_t) \propto \exp\left(-\eta \sum_{s=1}^t \ell_s(\theta)\right) \pi(\theta).$$

Not tractable in general, leading to variational approximations :

$$\begin{split} \tilde{\pi}_{t+1}(\theta) &= \operatorname*{arg\,min}_{q\in\mathcal{F}} \mathsf{KL}(q,\pi_{t+1}) \\ &= \operatorname*{arg\,min}_{q\in\mathcal{F}} \left\{ \mathbb{E}_{\theta\sim q} \left[\sum_{s=1}^t \ell_s(\theta) \right] + \frac{\mathsf{KL}(q,\pi)}{\eta} \right\}. \end{split}$$

Bayesian learning and variational inference (VI)

$$\pi_{t+1}(\theta) := \pi(\theta|x_1, y_1, \dots, x_t, y_t) \propto \exp\left(-\eta \sum_{s=1}^t \ell_s(\theta)\right) \pi(\theta).$$

Not tractable in general, leading to variational approximations :

$$\begin{split} \tilde{\pi}_{t+1}(\theta) &= \operatorname*{arg\,min}_{q\in\mathcal{F}} \mathsf{KL}(q,\pi_{t+1}) \\ &= \operatorname*{arg\,min}_{q\in\mathcal{F}} \left\{ \mathbb{E}_{\theta\sim q} \left[\sum_{s=1}^t \ell_s(\theta) \right] + \frac{\mathsf{KL}(q,\pi)}{\eta} \right\}. \end{split}$$

Formula for the online update of π_{t+1} :

$$\pi_{t+1}(\theta) \propto \exp\left(-\eta \ell_t(\theta)\right) \pi_t(\theta).$$

 $\mathbf{Q1}$: can we similarly define a sequential update for a variational approximation ?

Regret bounds for Bayesian inference

Theorem - classical, see Cesa-Bianchi & Lugosi 06

Under the assumption that the loss is bounded by B, the Bayesian update leads to

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t}[\ell_t(\theta)]$$

$$\leq \inf_q \left\{ \sum_{t=1}^{T} \mathbb{E}_{\theta \sim q}[\ell_t(\theta)] + \frac{\eta B^2 T}{8} + \frac{KL(q,\pi)}{\eta} \right\}.$$

N. Cesa-Bianchi & G. Lugosi (2006). Prediction, learning and games. Cambridge.

Explicit regret bounds

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t} [\ell_t(\theta)] \\ \leq \inf_{q} \left\{ \sum_{t=1}^{T} \mathbb{E}_{\theta \sim q} [\ell_t(\theta)] + \frac{\eta B^2 T}{8} + \frac{KL(q, \pi)}{\eta} \right\}.$$

Explicit regret bounds

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t} [\ell_t(\theta)]$$

$$\leq \inf_{q} \left\{ \sum_{t=1}^{T} \mathbb{E}_{\theta \sim q} [\ell_t(\theta)] + \frac{\eta B^2 T}{8} + \frac{KL(q, \pi)}{\eta} \right\}.$$

Example 1 : finite set $card(\Theta) = M$, prior π uniform. Consider $q = \delta_{\vartheta}$. Then :

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t}[\ell_t(\theta)] \leq \inf_{\vartheta} \sum_{t=1}^{T} \ell_t(\vartheta) + \frac{\eta B^2 T}{8} + \frac{\log(M)}{\eta}.$$

Explicit regret bounds

t=

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t} [\ell_t(\theta)] \\ \leq \inf_q \left\{ \sum_{t=1}^{T} \mathbb{E}_{\theta \sim q} [\ell_t(\theta)] + \frac{\eta B^2 T}{8} + \frac{\kappa L(q,\pi)}{\eta} \right\}.$$

Example 1 : finite set $card(\Theta) = M$, prior π uniform. Consider $q = \delta_{\vartheta}$. Then :

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t}[\ell_t(\theta)] \leq \inf_{\vartheta} \sum_{t=1}^{T} \ell_t(\vartheta) + \frac{\eta B^2 T}{8} + \frac{\log(M)}{\eta}$$

The choice $\eta = \sqrt{8 \log(M)/(TB^2)}$ leads to

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t}[\ell_t(\theta)] \leq \inf_{\vartheta} \sum_{t=1}^{T} \ell_t(\vartheta) + B\sqrt{\frac{T \log(M)}{2}}$$

Explicit regret bounds (2/2)

Example 2 : When $\Theta = \mathbb{R}^d$, using Gaussian priors and Gaussian q, and $\eta \sim \sqrt{T}$ "usually" leads to

$$\sum_{t=1}^T \mathbb{E}_{\theta \sim \pi_t}[\ell_t(\theta)] \leq \inf_{\|\theta\| \leq B} \sum_{t=1}^T \ell_t(\theta) + \mathcal{O}(\sqrt{dT\log(T)}).$$

Explicit regret bounds (2/2)

Example 2 : When $\Theta = \mathbb{R}^d$, using Gaussian priors and Gaussian q, and $\eta \sim \sqrt{T}$ "usually" leads to

$$\sum_{t=1}^T \mathbb{E}_{\theta \sim \pi_t}[\ell_t(\theta)] \leq \inf_{\|\theta\| \leq B} \sum_{t=1}^T \ell_t(\theta) + \mathcal{O}(\sqrt{dT\log(T)}).$$

$$\sum_{t=1}^{T} \mathbb{E}_{ heta \sim \pi_t}[\ell_t(heta)] \ \leq \inf_q \left\{ \sum_{t=1}^{T} \mathbb{E}_{ heta \sim q}[\ell_t(heta)] + rac{\eta B^2 T}{8} + rac{ extsf{KL}(q, \pi)}{\eta}
ight\}.$$

Q2 : can we derive similar results for online VI?

Parametric VI : $\mathcal{F} = \{q_{\mu}, \mu \in M\}.$

Parametric VI : $\mathcal{F} = \{q_{\mu}, \mu \in M\}.$

Sequential Variational Approximation (SVA) :

$$\theta_{t+1} = \arg\min_{\theta} \left\{ \left\langle \theta, \sum_{s=1}^{t} \nabla_{\theta} \ell_{s}(\theta_{s}) \right\rangle + \frac{\left\| \theta - \theta_{1} \right\|^{2}}{2\eta} \right\},\$$

Streaming Variational Bayes (SVB) :

$$heta_{t+1} = \operatorname*{arg\,min}_{ heta} \left\{ \left\langle heta,
abla_{ heta} \ell_t(heta_t) \right
angle + rac{\| heta - heta_t\|^2}{2\eta}
ight\},$$

Parametric VI : $\mathcal{F} = \{q_{\mu}, \mu \in M\}.$

Sequential Variational Approximation (SVA) :

$$\theta_{t+1} = \arg\min_{\theta} \left\{ \left\langle \theta, \sum_{s=1}^{t} \nabla_{\theta} \ell_{s}(\theta_{s}) \right\rangle + \frac{\|\theta - \theta_{1}\|^{2}}{2\eta} \right\},$$
$$\mu_{t+1} = \arg\min_{\mu} \left\{ \left\langle \mu, \sum_{s=1}^{t} \nabla_{\mu} \mathbb{E}_{\theta \sim q_{\mu_{s}}}[\ell_{s}(\theta)] \right\rangle + \frac{KL(q_{\mu}, \pi)}{\eta} \right\}.$$

Streaming Variational Bayes (SVB) :

$$heta_{t+1} = rgmin_{ heta} \left\{ \left\langle heta,
abla_{ heta} \ell_t(heta_t) \right
angle + rac{\| heta - heta_t\|^2}{2\eta}
ight\},$$

1

Parametric VI : $\mathcal{F} = \{q_{\mu}, \mu \in M\}.$

Sequential Variational Approximation (SVA) :

$$\theta_{t+1} = \arg\min_{\theta} \left\{ \left\langle \theta, \sum_{s=1}^{t} \nabla_{\theta} \ell_{s}(\theta_{s}) \right\rangle + \frac{\|\theta - \theta_{1}\|^{2}}{2\eta} \right\},$$
$$\mu_{t+1} = \arg\min_{\mu} \left\{ \left\langle \mu, \sum_{s=1}^{t} \nabla_{\mu} \mathbb{E}_{\theta \sim q_{\mu_{s}}}[\ell_{s}(\theta)] \right\rangle + \frac{KL(q_{\mu}, \pi)}{\eta} \right\}.$$

Streaming Variational Bayes (SVB) :

$$\theta_{t+1} = \arg\min_{\theta} \left\{ \left\langle \theta, \nabla_{\theta} \ell_t(\theta_t) \right\rangle + \frac{\|\theta - \theta_t\|^2}{2\eta} \right\},\$$
$$u_{t+1} = \arg\min_{\mu} \left\{ \left\langle \mu, \nabla_{\mu} \mathbb{E}_{\theta \sim q_{\mu_t}}[\ell_t(\theta)] \right\rangle + \frac{\mathsf{KL}(q_{\mu}, q_{\mu_t})}{\eta} \right\}.$$

SVA & SVB are tractable, and not equivalent

Example : Gaussian prior $\theta \sim \pi = \mathcal{N}(0, s^2 I)$ and mean-field Gaussian approximation, $\mu = (m, \sigma)$.

SVA : $m_{t+1} \leftarrow m_t - \eta s^2 \bar{g}_{m_t}$, $g_{t+1} \leftarrow g_t + \bar{g}_{\sigma_t}$, $\sigma_{t+1} \leftarrow h(\eta s g_{t+1}) s$, SVB : $m_{t+1} \leftarrow m_t - \eta \sigma_t^2 \bar{g}_{m_t}$, $\sigma_{t+1} \leftarrow \sigma_t h(\eta \sigma_t \bar{g}_{\sigma_t})$

where $h(x) := \sqrt{1 + x^2} - x$ is applied componentwise, as well as the multiplication of two vectors, and

$$\bar{g}_{m_t} = \frac{\partial}{\partial m} \mathbb{E}_{\theta \sim \pi_{m_t, \sigma_t}} [\ell_t(\theta)],$$
$$\bar{g}_{\sigma_t} = \frac{\partial}{\partial \sigma} \mathbb{E}_{\theta \sim \pi_{m_t, \sigma_t}} [\ell_t(\theta)].$$

Theoretical analysis of SVA

Theorem 1

Assume

•
$$\mu \mapsto \mathbb{E}_{\theta \sim q_{\mu}}[\ell_{t}(\theta)]$$
 is *L*-Lipschitz and convex,
• $\mu \mapsto \mathcal{K}L(q_{\mu}, \pi)$ is α -strongly convex, then

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim q_{\mu_{t}}}[\ell_{t}(\theta)]$$

$$\leq \inf_{\mu \in M} \left\{ \sum_{t=1}^{T} \mathbb{E}_{\theta \sim q_{\mu}}[\ell_{t}(\theta)] + \frac{\eta L^{2} T}{\alpha} + \frac{\mathcal{K}L(q_{\mu}, \pi)}{\eta} \right\}.$$

Theoretical analysis of SVA

Theorem 1

Assume

•
$$\mu \mapsto \mathbb{E}_{\theta \sim q_{\mu}}[\ell_{t}(\theta)]$$
 is *L*-Lipschitz and convex,
• $\mu \mapsto KL(q_{\mu}, \pi)$ is α -strongly convex, then

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim q_{\mu_{t}}}[\ell_{t}(\theta)]$$

$$\leq \inf_{\mu \in M} \left\{ \sum_{t=1}^{T} \mathbb{E}_{\theta \sim q_{\mu}}[\ell_{t}(\theta)] + \frac{\eta L^{2} T}{\alpha} + \frac{KL(q_{\mu}, \pi)}{\eta} \right\}.$$

Application to Gaussian approximation leads to

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim q_{\mu_t}}[\ell_t(\theta)] \leq \inf_{\theta} \sum_{t=1}^{T} \ell_t(\theta) + (1+o(1)) \frac{2L}{\alpha} \sqrt{dT \log(T)}.$$

The assumptions :

• $\mu \mapsto \mathbb{E}_{\theta \sim q_{\mu}}[\ell_t(\theta)]$ is *L*-Lipschitz and convex?

The assumptions :

• $\mu \mapsto \mathbb{E}_{\theta \sim q_{\mu}}[\ell_t(\theta)]$ is *L*-Lipschitz and convex?

Proposition

Assume $\theta \mapsto \ell_t(\theta)$ is L/2-Lipschitz and convex, and $\mu = (m, \Sigma)$ is a location scale parameter, then : satisfied.

The assumptions :

• $\mu \mapsto \mathbb{E}_{\theta \sim q_{\mu}}[\ell_t(\theta)]$ is *L*-Lipschitz and convex?

Proposition

Assume $\theta \mapsto \ell_t(\theta)$ is L/2-Lipschitz and convex, and $\mu = (m, \Sigma)$ is a location scale parameter, then : satisfied.

Proof : Lipschitz : in our paper ; convex :

J. Domke (2019). Provable smoothness guarantees for black-box variational inference.. Preprint arXiv.

The assumptions :

• $\mu \mapsto \mathbb{E}_{\theta \sim q_{\mu}}[\ell_t(\theta)]$ is *L*-Lipschitz and convex?

Proposition

Assume $\theta \mapsto \ell_t(\theta)$ is L/2-Lipschitz and convex, and $\mu = (m, \Sigma)$ is a location scale parameter, then : satisfied.

Proof : Lipschitz : in our paper; convex :

J. Domke (2019). Provable smoothness guarantees for black-box variational inference.. Preprint arXiv.

2)
$$\mu \mapsto \mathsf{KL}(q_{\mu}, \pi)$$
 is α -strongly convex?

The assumptions :

• $\mu \mapsto \mathbb{E}_{\theta \sim q_{\mu}}[\ell_t(\theta)]$ is *L*-Lipschitz and convex?

Proposition

Assume $\theta \mapsto \ell_t(\theta)$ is L/2-Lipschitz and convex, and $\mu = (m, \Sigma)$ is a location scale parameter, then : satisfied.

Proof : Lipschitz : in our paper; convex :

J. Domke (2019). Provable smoothness guarantees for black-box variational inference.. Preprint arXiv.

2
$$\mu \mapsto KL(q_{\mu}, \pi)$$
 is α -strongly convex?

 \rightarrow True for many examples, for example when q_{μ} and π are Gaussian (with lower and upper-bounded variance).

Theoretical analysis of SVB

Theorem 2

• We use Gaussian approximation q_{μ} with μ in a bounded set (diameter = D),

2 $\mu \mapsto \mathbb{E}_{\theta \sim q_{\mu}}[\ell_t(\theta)]$ is *L*-Lipschitz and convex,

3 $\mu \mapsto KL(q_{\mu}, \pi)$ is α -strongly convex, then

SVB with adequate η leads to

$$\sum_{t=1}^{T} \ell_t \Big(\mathbb{E}_{\theta \sim q_{\mu_t}}(\theta) \Big) \leq \inf_{\theta} \sum_{t=1}^{T} \ell_t(\theta) + DL\sqrt{2T}$$

Theoretical analysis of SVB

Theorem 2

• We use Gaussian approximation q_{μ} with μ in a bounded set (diameter = D),

2 $\mu \mapsto \mathbb{E}_{\theta \sim q_{\mu}}[\ell_t(\theta)]$ is *L*-Lipschitz and convex,

3 $\mu \mapsto KL(q_{\mu}, \pi)$ is α -strongly convex, then

SVB with adequate η leads to

$$\sum_{t=1}^{T} \ell_t \Big(\mathbb{E}_{\theta \sim q_{\mu_t}}(\theta) \Big) \leq \inf_{\theta} \sum_{t=1}^{T} \ell_t(\theta) + DL\sqrt{2T}$$

If, moreover, the loss is H-strongly convex,

$$\sum_{t=1}^{T} \ell_t \Big(\mathbb{E}_{\theta \sim q_{\mu_t}}(\theta) \Big) \leq \inf_{\theta} \sum_{t=1}^{T} \ell_t(\theta) + \frac{L^2(1 + \log(T))}{H}$$

Test on a simulated dataset



Figure – Average cumulative losses on different datasets for classification and regression tasks with OGA (yellow), OGA-EL (red), SVA (blue), SVB (purple) and NGVI (green).

Test on the Breast dataset



Figure – Average cumulative losses on different datasets for classification and regression tasks with OGA (yellow), OGA-EL (red), SVA (blue), SVB (purple) and NGVI (green).

Open questions



Analysis of SVB in the general case.

- Analysis of SVB in the general case.
- Analysis of the uncertainty quantification.

- Analysis of SVB in the general case.
- Analysis of the uncertainty quantification.
- In Solution NGVI is the next step in going closer to algorithms used to train Neural Networks with Bayesian principles. But it does not satisfy our assumptions...

NGVI

Assume we use approximations in the exponential family

$$q_{\lambda}(\theta) = a(\theta)b(\lambda) \exp\left[\langle \lambda, S(\theta) \rangle\right].$$
NGVI

Assume we use approximations in the exponential family

$$q_{\lambda}(\theta) = a(\theta)b(\lambda) \exp \left[\langle \lambda, S(\theta) \rangle\right].$$

We can re-parametrize this set with $\mu = \mathbb{E}_{\theta \sim q_{\lambda}}[T(\theta)] = F(\lambda)$.

NGVI

Assume we use approximations in the exponential family

$$q_{\lambda}(\theta) = a(\theta)b(\lambda) \exp \left[\langle \lambda, S(\theta) \rangle\right].$$

We can re-parametrize this set with $\mu = \mathbb{E}_{\theta \sim q_{\lambda}}[T(\theta)] = F(\lambda)$. The NGVI algorithm is given by

$$\lambda_{t+1} = (1-\rho)\lambda_t + \rho \nabla_{\mu} \mathbb{E}_{\theta \sim q_{\mu_t}} \left[\ell_t(\theta) \right],$$

M. E. Khan, D. Nielsen (2018). Fast yet Simple Natural-Gradient Descent for Variational Inference in Complex Models. ISITA.

NGVI

Assume we use approximations in the exponential family

$$q_{\lambda}(\theta) = a(\theta)b(\lambda) \exp\left[\langle \lambda, S(\theta) \rangle\right].$$

We can re-parametrize this set with $\mu = \mathbb{E}_{\theta \sim q_{\lambda}}[T(\theta)] = F(\lambda)$. The NGVI algorithm is given by

$$\lambda_{t+1} = (1-\rho)\lambda_t + \rho \nabla_{\mu} \mathbb{E}_{\theta \sim q_{\mu_t}} \left[\ell_t(\theta) \right],$$

M. E. Khan, D. Nielsen (2018). Fast yet Simple Natural-Gradient Descent for Variational Inference in Complex Models. ISITA.

Problem : under this parametrisation, $\mu \mapsto \mathbb{E}_{\theta \sim q_{\lambda}}[\ell_t(\theta)]$ is generally not convex...

Thank you!