

# Generalization bounds for variational inference

Pierre Alquier



3rd International Conference on  
Econometrics and Statistics (EcoSta 2019)  
25-27 June 2019, National Chung Hsing University, Taiwan



# Notations

Assume that we observe  $X_1, \dots, X_n$  i.i.d from  $P_{\theta_0}$  in a model  $\{P_{\theta}, \theta \in \Theta\}$  dominated by  $Q : \frac{dP_{\theta}}{dQ} = p_{\theta}$ . Prior  $\pi$  on  $\Theta$ .

# Notations

Assume that we observe  $X_1, \dots, X_n$  i.i.d from  $P_{\theta_0}$  in a model  $\{P_{\theta}, \theta \in \Theta\}$  dominated by  $Q : \frac{dP_{\theta}}{dQ} = p_{\theta}$ . Prior  $\pi$  on  $\Theta$ .

## The likelihood

$$L_n(\theta) = \prod_{i=1}^n p_{\theta}(X_i)$$

# Notations

Assume that we observe  $X_1, \dots, X_n$  i.i.d from  $P_{\theta_0}$  in a model  $\{P_{\theta}, \theta \in \Theta\}$  dominated by  $Q : \frac{dP_{\theta}}{dQ} = p_{\theta}$ . Prior  $\pi$  on  $\Theta$ .

## The likelihood

$$L_n(\theta) = \prod_{i=1}^n p_{\theta}(X_i)$$

## The posterior

$$\pi_n(d\theta) \propto L_n(\theta)\pi(d\theta).$$

# Notations

Assume that we observe  $X_1, \dots, X_n$  i.i.d from  $P_{\theta_0}$  in a model  $\{P_{\theta}, \theta \in \Theta\}$  dominated by  $Q : \frac{dP_{\theta}}{dQ} = p_{\theta}$ . Prior  $\pi$  on  $\Theta$ .

## The likelihood

$$L_n(\theta) = \prod_{i=1}^n p_{\theta}(X_i)$$

## The posterior

$$\pi_n(d\theta) \propto L_n(\theta)\pi(d\theta).$$

## The tempered posterior - $0 < \alpha \leq 1$

$$\pi_{n,\alpha}(d\theta) \propto [L_n(\theta)]^{\alpha}\pi(d\theta).$$

# Computation of the posterior

- explicit form (conjugate models),

# Computation of the posterior

- explicit form (conjugate models),
- MCMC algorithms : Metropolis-Hastings, Gibbs sampler, Langevin Monte Carlo...

# Computation of the posterior

- explicit form (conjugate models),
- MCMC algorithms : Metropolis-Hastings, Gibbs sampler, Langevin Monte Carlo...

But...

- when the dimension is large, the convergence of MCMC can be extremely slow,

# Computation of the posterior

- explicit form (conjugate models),
- MCMC algorithms : Metropolis-Hastings, Gibbs sampler, Langevin Monte Carlo...

But...

- when the dimension is large, the convergence of MCMC can be extremely slow,
- when the model is complex or when the sample size is large, each evaluation of  $\pi_{n,\alpha}(\theta)$  can be expensive.

# Computation of the posterior

- explicit form (conjugate models),
- MCMC algorithms : Metropolis-Hastings, Gibbs sampler, Langevin Monte Carlo...

But...

- when the dimension is large, the convergence of MCMC can be extremely slow,
- when the model is complex or when the sample size is large, each evaluation of  $\pi_{n,\alpha}(\theta)$  can be expensive.

For these reasons, in the past 20 years, many methods targeting an approximation of  $\pi_{n,\alpha}$  became popular : ABC, EP algorithm, **variational inference**, approximate MCMC ...

# Variational approximations : definitions

**Idea of VB** : chose a family  $\mathcal{F}$  of probability distributions on  $\Theta$  and approximate  $\pi_{n,\alpha}$  by a distribution in  $\mathcal{F}$  :

# Variational approximations : definitions

**Idea of VB** : chose a family  $\mathcal{F}$  of probability distributions on  $\Theta$  and approximate  $\pi_{n,\alpha}$  by a distribution in  $\mathcal{F}$  :

$$\tilde{\pi}_{n,\alpha} := \arg \min_{\rho \in \mathcal{F}} \mathcal{K}(\rho, \pi_{n,\alpha}).$$

# Variational approximations : definitions

**Idea of VB** : chose a family  $\mathcal{F}$  of probability distributions on  $\Theta$  and approximate  $\pi_{n,\alpha}$  by a distribution in  $\mathcal{F}$  :

$$\tilde{\pi}_{n,\alpha} := \arg \min_{\rho \in \mathcal{F}} \mathcal{K}(\rho, \pi_{n,\alpha}).$$

Examples :

- parametric approximation

$$\mathcal{F} = \{ \mathcal{N}(\mu, \Sigma) : \mu \in \mathbb{R}^d, \Sigma \in \mathcal{S}_d^+ \}.$$

# Variational approximations : definitions

**Idea of VB :** chose a family  $\mathcal{F}$  of probability distributions on  $\Theta$  and approximate  $\pi_{n,\alpha}$  by a distribution in  $\mathcal{F}$  :

$$\tilde{\pi}_{n,\alpha} := \arg \min_{\rho \in \mathcal{F}} \mathcal{K}(\rho, \pi_{n,\alpha}).$$

Examples :

- parametric approximation

$$\mathcal{F} = \{ \mathcal{N}(\mu, \Sigma) : \mu \in \mathbb{R}^d, \Sigma \in \mathcal{S}_d^+ \}.$$

- mean-field approximation,  $\Theta = \Theta_1 \times \Theta_2$  and

$$\mathcal{F} : \{ \rho : \rho(d\theta) = \rho_1(d\theta_1) \times \rho_2(d\theta_2) \}.$$

# Empirical lower bound (ELBO)

Note that :

$$\begin{aligned}\tilde{\pi}_{n,\alpha} &= \arg \min_{\rho \in \mathcal{F}} \mathcal{K}(\rho, \pi_{n,\alpha}) \\ &= \arg \min_{\rho \in \mathcal{F}} \underbrace{\left\{ -\alpha \int \frac{1}{n} \sum_{i=1}^n \log p_{\theta}(X_i) \rho(d\theta) + \mathcal{K}(\rho, \pi) \right\}}_{-\text{ELBO}(\rho)}.\end{aligned}$$

# Empirical lower bound (ELBO)

Note that :

$$\begin{aligned}\tilde{\pi}_{n,\alpha} &= \arg \min_{\rho \in \mathcal{F}} \mathcal{K}(\rho, \pi_{n,\alpha}) \\ &= \arg \min_{\rho \in \mathcal{F}} \underbrace{\left\{ -\alpha \int \frac{1}{n} \sum_{i=1}^n \log p_{\theta}(X_i) \rho(d\theta) + \mathcal{K}(\rho, \pi) \right\}}_{-\text{ELBO}(\rho)}.\end{aligned}$$

So we have the equivalent definition :

$$\tilde{\pi}_{n,\alpha} := \arg \max_{\rho \in \mathcal{F}} \text{ELBO}(\rho).$$

# Consistency results

3 papers (2017) :

- 1 consistency and rates of convergence for  $\alpha < 1$  :



P. Alquier & J. Ridgway. Concentration of tempered posteriors and of their variational approximations. To appear in *The Annals of Statistics*.

# Consistency results

3 papers (2017) :

- 1 consistency and rates of convergence for  $\alpha < 1$  :



P. Alquier & J. Ridgway. Concentration of tempered posteriors and of their variational approximations. To appear in *The Annals of Statistics*.

- 2 extension to models with hidden variables :



A. Bhattacharya, D. Pati & Y. Yang. Bayesian fractional posteriors. *The Annals of Statistics*, 2019.

# Consistency results

3 papers (2017) :

① consistency and rates of convergence for  $\alpha < 1$  :



P. Alquier & J. Ridgway. Concentration of tempered posteriors and of their variational approximations. To appear in *The Annals of Statistics*.

② extension to models with hidden variables :



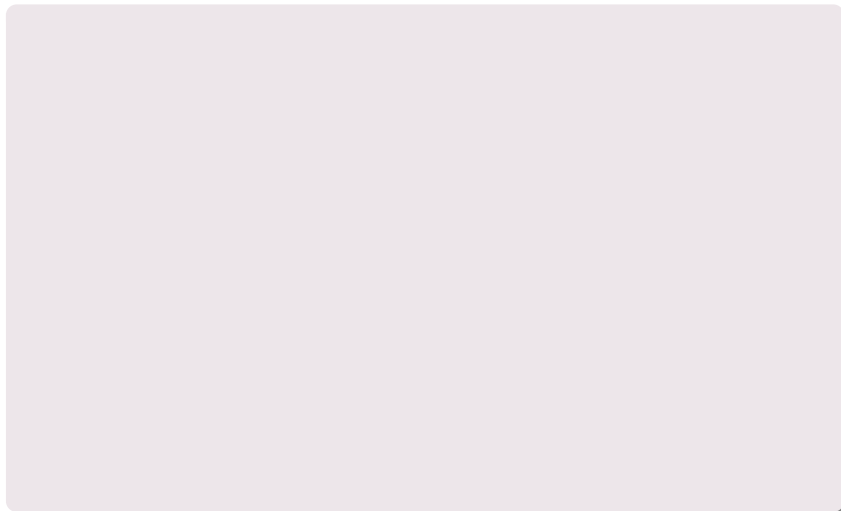
A. Bhattacharya, D. Pati & Y. Yang. Bayesian fractional posteriors. *The Annals of Statistics*, 2019.

③ extension to  $\alpha = 1$  :



F. Zhang & C. Gao. Convergence Rates of Variational Posterior Distributions. *Preprint arXiv*, 2017.

# Sequential estimation problem



# Sequential estimation problem

- 1 initialize  $\theta_1$ ,

# Sequential estimation problem

- 1 initialize  $\theta_1$ ,
- 2  $x_1$  revealed,

# Sequential estimation problem

- 1
- 1 initialize  $\theta_1$ ,
- 2  $x_1$  revealed,
- 3 incur loss  
     $-\log p_{\theta_1}(x_1)$

# Sequential estimation problem

- 1
  - 1 initialize  $\theta_1$ ,
  - 2  $x_1$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_1}(x_1)$
- 2
  - 1 update  $\theta_1 \rightarrow \theta_2$ ,

# Sequential estimation problem

- 1
  - 1 initialize  $\theta_1$ ,
  - 2  $x_1$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_1}(x_1)$
- 2
  - 1 update  $\theta_1 \rightarrow \theta_2$ ,
  - 2  $x_2$  revealed,

# Sequential estimation problem

- 1
  - 1 initialize  $\theta_1$ ,
  - 2  $x_1$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_1}(x_1)$
- 2
  - 1 update  $\theta_1 \rightarrow \theta_2$ ,
  - 2  $x_2$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_2}(x_2)$

# Sequential estimation problem

- 1
  - 1 initialize  $\theta_1$ ,
  - 2  $x_1$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_1}(x_1)$
- 2
  - 1 update  $\theta_1 \rightarrow \theta_2$ ,
  - 2  $x_2$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_2}(x_2)$
- 3
  - 1 update  $\theta_2 \rightarrow \theta_3$ ,

# Sequential estimation problem

- 1
  - 1 initialize  $\theta_1$ ,
  - 2  $x_1$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_1}(x_1)$
- 2
  - 1 update  $\theta_1 \rightarrow \theta_2$ ,
  - 2  $x_2$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_2}(x_2)$
- 3
  - 1 update  $\theta_2 \rightarrow \theta_3$ ,
  - 2  $x_3$  revealed,

# Sequential estimation problem

- 1
  - 1 initialize  $\theta_1$ ,
  - 2  $x_1$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_1}(x_1)$
- 2
  - 1 update  $\theta_1 \rightarrow \theta_2$ ,
  - 2  $x_2$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_2}(x_2)$
- 3
  - 1 update  $\theta_2 \rightarrow \theta_3$ ,
  - 2  $x_3$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_3}(x_3)$
- 4 ...

# Sequential estimation problem

Objective :

- 1
  - 1 initialize  $\theta_1$ ,
  - 2  $x_1$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_1}(x_1)$
- 2
  - 1 update  $\theta_1 \rightarrow \theta_2$ ,
  - 2  $x_2$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_2}(x_2)$
- 3
  - 1 update  $\theta_2 \rightarrow \theta_3$ ,
  - 2  $x_3$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_3}(x_3)$
- 4 ...

# Sequential estimation problem

- 1
  - 1 initialize  $\theta_1$ ,
  - 2  $x_1$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_1}(x_1)$
- 2
  - 1 update  $\theta_1 \rightarrow \theta_2$ ,
  - 2  $x_2$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_2}(x_2)$
- 3
  - 1 update  $\theta_2 \rightarrow \theta_3$ ,
  - 2  $x_3$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_3}(x_3)$

4 ...

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

# Sequential estimation problem

- 1
  - 1 initialize  $\theta_1$ ,
  - 2  $x_1$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_1}(x_1)$
- 2
  - 1 update  $\theta_1 \rightarrow \theta_2$ ,
  - 2  $x_2$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_2}(x_2)$
- 3
  - 1 update  $\theta_2 \rightarrow \theta_3$ ,
  - 2  $x_3$  revealed,
  - 3 incur loss  
 $-\log p_{\theta_3}(x_3)$

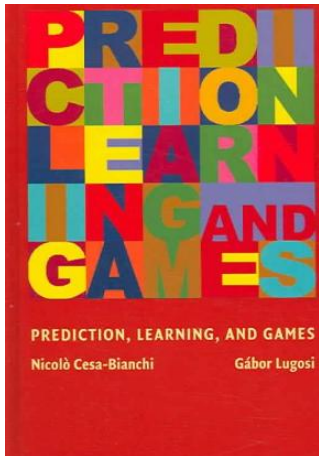
4 ...

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

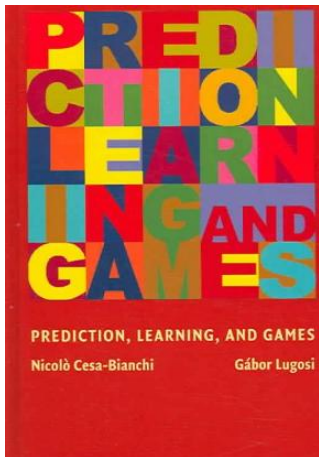
$$\sum_{t=1}^T [-\log p_{\theta_t}(x_t)]$$

as small as possible for any  $T$ ,  
**without stochastic**  
**assumptions on the data.**

# Reference



# Reference



The regret :

$$R(T) = \sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \\ - \inf_{\theta \in \Theta} \sum_{t=1}^T [-\log p_{\theta}(x_t)].$$

# EWA strategy / multiplicative update...

# EWA strategy / multiplicative update...

- learning rate  $\alpha > 0$ .

# EWA strategy / multiplicative update...

- learning rate  $\alpha > 0$ .
- initialize  $p_1 = \pi$  (the prior).

# EWA strategy / multiplicative update...

- learning rate  $\alpha > 0$ .
- initialize  $p_1 = \pi$  (the prior).

---

## Algorithm 2 Exponentially Weighted Aggregation

---

- 1: **for**  $t = 1, 2, \dots$  **do**
  - 2:    $\theta_t = \mathbb{E}_{\theta \sim p_t}[\theta]$ ,
  - 3:    $x_t$  revealed, update  $p_{t+1}(d\theta) = \frac{[p_{\theta}(x_t)]^{\alpha} p_t(d\theta)}{\int [p_{\vartheta}(x_t)]^{\alpha} p_t(d\vartheta)}$ .
  - 4: **end for**
-

# EWA strategy / multiplicative update...

- learning rate  $\alpha > 0$ .
- initialize  $p_1 = \pi$  (the prior).

---

## Algorithm 2 Exponentially Weighted Aggregation

---

- 1: **for**  $t = 1, 2, \dots$  **do**
  - 2:    $\theta_t = \mathbb{E}_{\theta \sim p_t}[\theta]$ ,
  - 3:    $x_t$  revealed, update  $p_{t+1}(d\theta) = \frac{[p_\theta(x_t)]^\alpha p_t(d\theta)}{\int [p_{\vartheta}(x_t)]^\alpha p_t(d\vartheta)}$ .
  - 4: **end for**
- 

Note that  $p_t = \pi_{n,\alpha}$  the tempered posterior, so problem : how can we compute  $\theta_t$ ?

# A regret bound for EWA

From now,  $\theta \mapsto [-\log p_\theta(x_t)]$  is convex + bounded :  $|\cdot| \leq C$ .

# A regret bound for EWA

From now,  $\theta \mapsto [-\log p_\theta(x_t)]$  is convex + bounded :  $|\cdot| \leq C$ .

## Theorem

$$\sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \leq \inf_p \left[ \sum_{t=1}^T \mathbb{E}_{\theta \sim p} [-\log p_\theta(x_t)] + \frac{\alpha C^2 T}{2} + \frac{\mathcal{K}(p, \pi)}{\alpha} \right].$$

# A regret bound for EWA

From now,  $\theta \mapsto [-\log p_\theta(x_t)]$  is convex + bounded :  $|\cdot| \leq C$ .

## Theorem

$$\sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \leq \inf_p \left[ \sum_{t=1}^T \mathbb{E}_{\theta \sim p} [-\log p_\theta(x_t)] + \frac{\alpha C^2 T}{2} + \frac{\mathcal{K}(p, \pi)}{\alpha} \right].$$

Under similar assumptions than in the batch case, that is, the prior gives enough mass to relevant  $\theta$ , and  $\alpha \sim 1/\sqrt{T}$ ,

$$\sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \leq \inf_{\theta \in \Theta} \sum_{t=1}^T [-\log p_\theta(x_t)] + \text{cst.} \sqrt{T}$$

$$\sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \leq \inf_{\theta \in \Theta} \sum_{t=1}^T [-\log p_{\theta}(x_t)] + \text{cst.} \sqrt{T}$$

$$\sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \leq \inf_{\theta \in \Theta} \sum_{t=1}^T [-\log p_{\theta}(x_t)] + \text{cst} \cdot \sqrt{T}$$

$$\frac{1}{T} \sum_{t=1}^T \log \frac{q(x_t)}{p_{\theta_t}(x_t)} \leq \inf_{\theta \in \Theta} \frac{1}{T} \sum_{t=1}^T \log \frac{q(x_t)}{p_{\theta}(x_t)} + \frac{\text{cst}}{\sqrt{T}}.$$

$$\sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \leq \inf_{\theta \in \Theta} \sum_{t=1}^T [-\log p_{\theta}(x_t)] + \text{cst} \cdot \sqrt{T}$$

$$\frac{1}{T} \sum_{t=1}^T \log \frac{q(x_t)}{p_{\theta_t}(x_t)} \leq \inf_{\theta \in \Theta} \frac{1}{T} \sum_{t=1}^T \log \frac{q(x_t)}{p_{\theta}(x_t)} + \frac{\text{cst}}{\sqrt{T}}.$$

Assuming that  $x_1, \dots, x_T$  are actually i.i.d from  $Q$ , with density  $q$ , define

$$\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \theta_t,$$

we have (“online-to-batch” conversion) :

$$\mathbb{E} [\mathcal{K}(Q, P_{\hat{\theta}_T})] \leq \inf_{\theta \in \Theta} \mathcal{K}(Q, P_{\theta}) + \frac{\text{cst}}{\sqrt{T}}.$$

# Variational approximations of EWA



B.-E. Chérif-Abdellatif, P. Alquier & M. E. Khan. A Generalization Bound for Online Variational Inference. *Preprint arXiv*, 2019.

# Variational approximations of EWA



B.-E. Chérif-Abdellatif, P. Alquier & M. E. Khan. A Generalization Bound for Online Variational Inference. *Preprint arXiv*, 2019.



# Variational approximations of EWA



B.-E. Chérif-Abdellatif, P. Alquier & M. E. Khan. A Generalization Bound for Online Variational Inference. *Preprint arXiv*, 2019.



Parametric variational approximation :  $\mathcal{F} = \{q_\mu, \mu \in M\}$ .  
Objective : propose a way to update  $\mu_t \rightarrow \mu_{t+1}$  so that  $q_{\mu_t}$  leads to similar performances as  $p_t$  in EWA...

# SVA and SVB strategies

---

**Algorithm 3** SVA (Sequential Variational Approximation)

---

- 1: **for**  $t = 1, 2, \dots$  **do**
- 2:    $\theta_t = \mathbb{E}_{\theta \sim q_{\mu_t}}[\theta],$
- 3:    $x_t$  revealed, update

$$\mu_{t+1} = \arg \min_{\mu \in M} \left[ \mu^T \nabla_{\mu} \sum_{i=1}^t \mathbb{E}_{\theta \sim q_{\mu}} [-\log p_{\theta}(x_i)] + \frac{\mathcal{K}(q_{\mu}, \pi)}{\alpha} \right].$$

- 4: **end for**
-

# SVA and SVB strategies

---

**Algorithm 3** SVA (Sequential Variational Approximation)

---

- 1: **for**  $t = 1, 2, \dots$  **do**
- 2:    $\theta_t = \mathbb{E}_{\theta \sim q_{\mu_t}}[\theta],$
- 3:    $x_t$  revealed, update

$$\mu_{t+1} = \arg \min_{\mu \in M} \left[ \mu^T \nabla_{\mu} \sum_{i=1}^t \mathbb{E}_{\theta \sim q_{\mu}} [-\log p_{\theta}(x_i)] + \frac{\mathcal{K}(q_{\mu}, \pi)}{\alpha} \right].$$

- 4: **end for**
- 

SVB (Streaming Variational Bayes) has update

$$\mu_{t+1} = \arg \min_{\mu \in M} \left[ \mu^T \nabla_{\mu} \mathbb{E}_{\theta \sim q_{\mu}} [-\log p_{\theta}(x_t)] + \frac{\mathcal{K}(q_{\mu}, q_{\mu_t})}{\alpha} \right].$$

# NGVI strategy

NGVI (Natural Gradient Variational Inference) : fix some  $\beta > 0$ ,

$$\mu_{t+1} = \arg \min_{\mu \in M} \left[ \mu^T \nabla_{\mu} \mathbb{E}_{\theta \sim q_{\mu}} [-\log p_{\theta}(x_t)] + \frac{\mathcal{K}(q_{\mu}, \pi)}{\alpha} + \frac{\mathcal{K}(q_{\mu}, q_{\mu_t})}{\beta} \right].$$

# NGVI strategy

NGVI (Natural Gradient Variational Inference) : fix some  $\beta > 0$ ,

$$\mu_{t+1} = \arg \min_{\mu \in M} \left[ \mu^T \nabla_{\mu} \mathbb{E}_{\theta \sim q_{\mu}} [-\log p_{\theta}(x_t)] + \frac{\mathcal{K}(q_{\mu}, \pi)}{\alpha} + \frac{\mathcal{K}(q_{\mu}, q_{\mu_t})}{\beta} \right].$$



M. E. Khan & W. Lin. Conjugate-computation variational inference : Converting variational inference in non-conjugate models to inferences in conjugate models. *AISTAT*, 2017.

# An example : SVB with Gaussian approximations

As an example, assume that  $\theta \in \mathbb{R}^d$ , the prior is  $\pi = \mathcal{N}(0, s^2 I)$  and that we use the variational approximation

$$\text{family : } q_{\mu} = q_{m, \sigma} = \mathcal{N} \left( m, \begin{pmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_d^2 \end{pmatrix} \right).$$

# An example : SVB with Gaussian approximations

As an example, assume that  $\theta \in \mathbb{R}^d$ , the prior is  $\pi = \mathcal{N}(0, s^2 I)$  and that we use the variational approximation

$$\text{family : } q_{\mu} = q_{m, \sigma} = \mathcal{N} \left( m, \begin{pmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_d^2 \end{pmatrix} \right).$$

In this case, the update in SVB is :

$$\begin{aligned} m_{t+1} &= m_t - \alpha \sigma_t^2 \odot \nabla_{m=m_t} \mathbb{E}_{\theta \sim q_{m, \sigma_t}} [-\log p_{\theta}(x_t)] \\ \sigma_{t+1} &= \sigma_t \odot h \left( \frac{\alpha \sigma_t \nabla_{\sigma=\sigma_t} \mathbb{E}_{\theta \sim q_{m_t, \sigma}} [-\log p_{\theta}(x_t)]}{2} \right) \end{aligned}$$

where  $\odot$  means “componentwise multiplication” and  $h(x) = \sqrt{1 + x^2} - x$  is also applied componentwise.

# An example : SVB with Gaussian approximations

As an example, assume that  $\theta \in \mathbb{R}^d$ , the prior is  $\pi = \mathcal{N}(0, s^2 I)$  and that we use the variational approximation

$$\text{family : } q_{\mu} = q_{m, \sigma} = \mathcal{N} \left( m, \begin{pmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_d^2 \end{pmatrix} \right).$$

In this case, the update in SVB is :

$$m_{t+1} = m_t - \alpha \sigma_t^2 \odot \nabla_{m=m_t} \mathbb{E}_{\theta \sim q_{m_t, \sigma_t}} [-\log p_{\theta}(x_t)]$$
$$\sigma_{t+1} = \sigma_t \odot h \left( \frac{\alpha \sigma_t \nabla_{\sigma=\sigma_t} \mathbb{E}_{\theta \sim q_{m_t, \sigma}} [-\log p_{\theta}(x_t)]}{2} \right)$$

where  $\odot$  means “componentwise multiplication” and  $h(x) = \sqrt{1+x^2} - x$  is also applied componentwise. We also have explicit formulas for SVA and NGVI (see the paper).

# A regret bound for SVA

## Theorem (Chérif-Abdellatif, A. & Khan)

Assume that  $\mu \mapsto \mathbb{E}_{\theta \sim q_\mu} [-\log p_\theta(x_t)]$  is  $L$ -Lipschitz and convex.

# A regret bound for SVA

## Theorem (Chérif-Abdellatif, A. & Khan)

Assume that  $\mu \mapsto \mathbb{E}_{\theta \sim q_\mu} [-\log p_\theta(x_t)]$  is  $L$ -Lipschitz and convex. (this is for example the case as soon as the log-likelihood is concave in  $\theta$  and  $L$ -Lipschitz, and  $\mu$  is a location-scale parameter).

# A regret bound for SVA

## Theorem (Chérif-Abdellatif, A. & Khan)

Assume that  $\mu \mapsto \mathbb{E}_{\theta \sim q_\mu} [-\log p_\theta(x_t)]$  is  $L$ -Lipschitz and convex. Assume that  $\mu \mapsto \mathcal{K}(p_\mu, \pi)$  is  $\gamma$ -strongly convex. Then SVA satisfies :

$$\begin{aligned} & \sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \\ & \leq \inf_{\mu \in M} \left\{ \mathbb{E}_{\theta \sim q_\mu} \left[ \sum_{t=1}^T [-\log p_\theta(x_t)] \right] + \frac{\alpha L^2 T}{\gamma} + \frac{\mathcal{K}(q_\mu, \pi)}{\alpha} \right\}. \end{aligned}$$

# A regret bound for SVA

## Theorem (Chérif-Abdellatif, A. & Khan)

Assume that  $\mu \mapsto \mathbb{E}_{\theta \sim q_\mu} [-\log p_\theta(x_t)]$  is  $L$ -Lipschitz and convex. Assume that  $\mu \mapsto \mathcal{K}(p_\mu, \pi)$  is  $\gamma$ -strongly convex. Then SVA satisfies :

$$\begin{aligned} & \sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \\ & \leq \inf_{\mu \in M} \left\{ \mathbb{E}_{\theta \sim q_\mu} \left[ \sum_{t=1}^T [-\log p_\theta(x_t)] \right] + \frac{\alpha L^2 T}{\gamma} + \frac{\mathcal{K}(q_\mu, \pi)}{\alpha} \right\}. \end{aligned}$$

For SVB : some results in the Gaussian case.

# A regret bound for SVA

## Theorem (Chérif-Abdellatif, A. & Khan)

Assume that  $\mu \mapsto \mathbb{E}_{\theta \sim q_\mu} [-\log p_\theta(x_t)]$  is  $L$ -Lipschitz and convex. Assume that  $\mu \mapsto \mathcal{K}(p_\mu, \pi)$  is  $\gamma$ -strongly convex. Then SVA satisfies :

$$\begin{aligned} & \sum_{t=1}^T [-\log p_{\theta_t}(x_t)] \\ & \leq \inf_{\mu \in M} \left\{ \mathbb{E}_{\theta \sim q_\mu} \left[ \sum_{t=1}^T [-\log p_\theta(x_t)] \right] + \frac{\alpha L^2 T}{\gamma} + \frac{\mathcal{K}(q_\mu, \pi)}{\alpha} \right\}. \end{aligned}$$

For SVB : some results in the Gaussian case. For NGVI : we were not able to derive regret bounds until now.

# Test on a simulated dataset

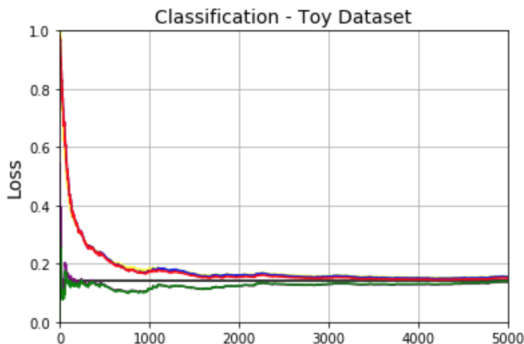


Figure – Average cumulative losses on different datasets for classification and regression tasks with OGA (yellow), OGA-EL (red), SVA (blue), SVB (purple) and NGVI (green).

# Test on the Breast dataset

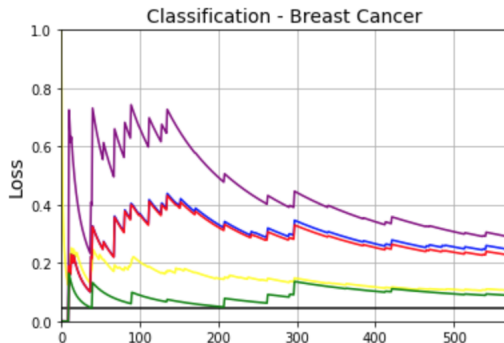


Figure – Average cumulative losses on different datasets for classification and regression tasks with OGA (yellow), OGA-EL (red), SVA (blue), SVB (purple) and NGVI (green).

# Test on the Pima Indians dataset

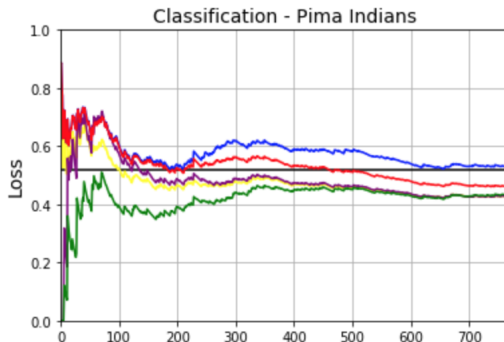


Figure – Average cumulative losses on different datasets for classification and regression tasks with OGA (yellow), OGA-EL (red), SVA (blue), SVB (purple) and NGVI (green).

# Test on the Boston Housing dataset

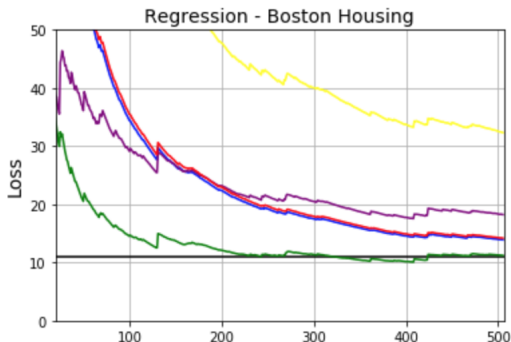
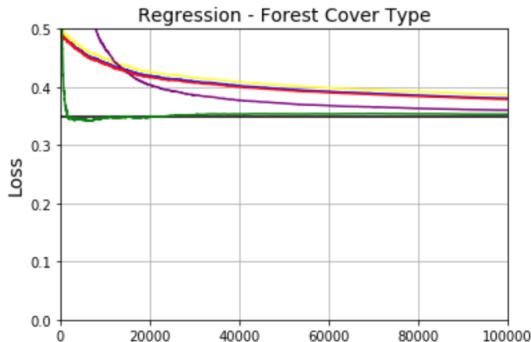


Figure – Average cumulative losses on different datasets for classification and regression tasks with OGA (yellow), OGA-EL (red), SVA (blue), SVB (purple) and NGVI (green).

# Test on the Forest Cover Type 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).

# Conclusions

- 1 Using online-to-batch conversion, we now have algorithms for variational inference with provable statistical properties after a finite number of steps.

# Conclusions

- 1 Using online-to-batch conversion, we now have algorithms for variational inference with provable statistical properties after a finite number of steps.
- 2 SVA, SVB competitive with OGA (online gradient algorithm, “non-Bayesian”).

# Conclusions

- ➊ Using online-to-batch conversion, we now have algorithms for variational inference with provable statistical properties after a finite number of steps.
- ➋ SVA, SVB competitive with OGA (online gradient algorithm, “non-Bayesian”).
- ➌ NGVI is the best method on all datasets. Its theoretical analysis is thus an important open problem. Cannot be done with our current techniques (using natural parameters in exponential models lead to non-convex objectives).

Thank you !