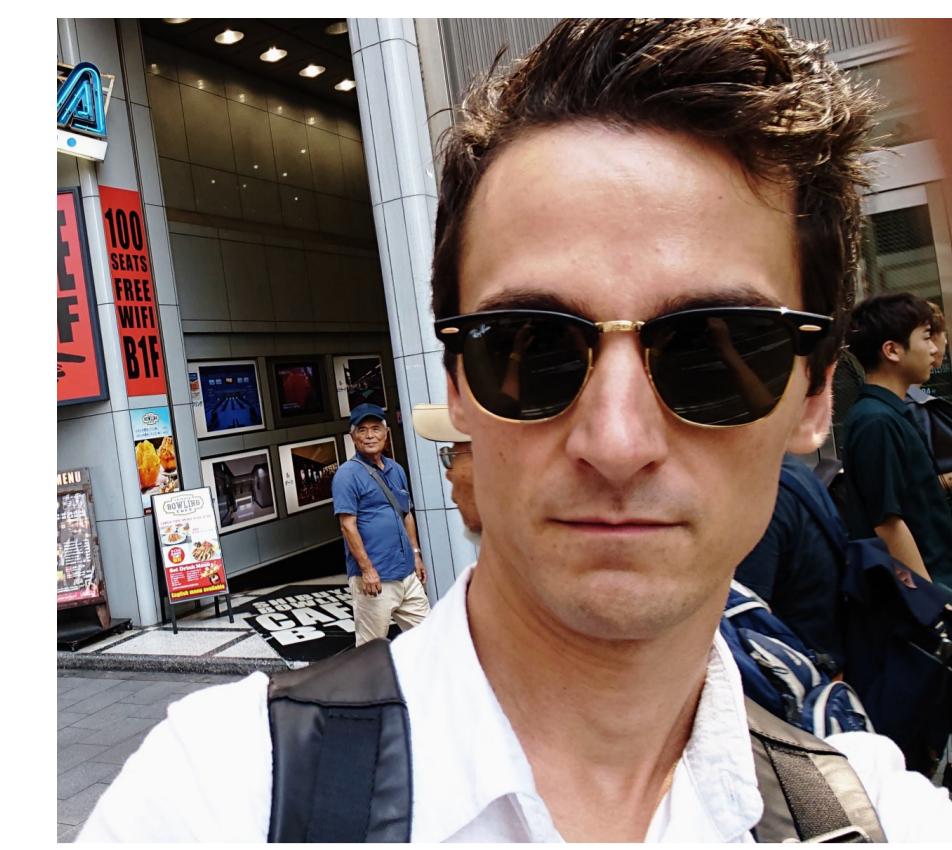


# Non-Exponentially Weighted Aggregation: Regret Bounds for Unbounded Loss Functions

Pierre Alquier  
RIKEN AIP



## (Generalized) Bayes update

$$\rho^t = \arg \min_{\rho} \left\{ \sum_{s=1}^{t-1} \mathbb{E}_{\theta \sim \rho} [\ell_s(\theta)] + \frac{\text{KL}(\rho \parallel \pi)}{\eta} \right\}.$$

► no constraint on  $\rho$ :

$$\rho^t(d\theta) \propto \exp \left[ -\eta \sum_{s=1}^{t-1} \ell_s(\theta) \right] \pi(d\theta).$$

► constraint on  $\rho$ : variational inference.

Reasons to go beyond KL:

KNOBLAUCH, J., JEWSON, J. & DAMOULAS, T. (2019). Generalized variational inference: Three arguments for deriving new posteriors. *Preprint arXiv*.

## Objective

$$\rho^t = \arg \min_{\rho} \left\{ \sum_{s=1}^{t-1} \mathbb{E}_{\theta \sim \rho} [\ell_s(\theta)] + \frac{D_{\phi}(\rho \parallel \pi)}{\eta} \right\}.$$

$$D_{\phi}(\rho \parallel \pi) = \mathbb{E}_{\theta \sim \pi} \left[ \phi \left( \frac{d\rho}{d\pi}(\theta) \right) \right]$$

► formula for the update?

► regret bounds?

## Theorem: formula for $\rho^t$ – “non-exponential weights”

$$\nabla \tilde{\phi}^*(y) = \arg \max_{x \geq 0} \{xy - \phi(x)\},$$

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

The proof uses convex analysis tools from:

AGRAWAL, R. & HOREL, T. (2020). Optimal Bounds between  $f$ -Divergences and Integral Probability Metrics. *ICML*.

## Example 1: $D_{\phi}(\rho \parallel \pi) = \text{KL}(\rho \parallel \pi)$

$$\phi(x) = x \log(x)$$

$$\nabla \tilde{\phi}^*(y) = \exp(y) - 1$$

$$\rho^t(d\theta) = \exp \left[ \lambda_t - \eta \sum_{s=1}^{t-1} \ell_s(\theta) - 1 \right] \pi(d\theta)$$

## Example 2: $D_{\phi}(\rho \parallel \pi) = \chi^2(\rho \parallel \pi)$

$$\begin{aligned} \Phi(x) &= x^2 - 1 \\ \nabla \tilde{\phi}^*(y) &= \max(0, y/2) \end{aligned}$$

$$\rho^t(d\theta) = \frac{1}{2} \max \left[ 0, \lambda_t - \eta \sum_{s=1}^{t-1} \ell_s(\theta) \right] \pi(d\theta)$$

### Theorem: regret bound

Assume there is a norm  $\|\cdot\|$  such that

1.  $\rho \mapsto \mathbb{E}_{\theta \sim \rho} [\ell_t(\theta)]$  is  $L$ -Lipschitz w.r.t  $\|\cdot\|$ ,
2.  $\rho \mapsto D_{\phi}(\rho \parallel \pi)$  is  $\alpha$ -strongly convex w.r.t  $\|\cdot\|$ .

$$\sum_{t=1}^T \mathbb{E}_{\theta \sim \rho^t} [\ell_t(\theta)] \leq \inf_{\rho \in \mathcal{P}(\Theta)} \left\{ \sum_{t=1}^T \mathbb{E}_{\theta \sim \rho} [\ell_t(\theta)] + \frac{\eta L^2 T}{\alpha} + \frac{D_{\phi}(\rho \parallel \pi)}{\eta} \right\}.$$

## Example 1: $D_{\phi}(\rho \parallel \pi) = \text{KL}(\rho \parallel \pi)$

► known result:  $\text{KL}(\rho \parallel \pi)$  is 1-strongly convex with respect to  $\|\cdot\|_{\text{TV}}$ ;

► we have:

$$\begin{aligned} \left| \int \ell_t(\theta) \rho(d\theta) - \int \ell_t \rho'(d\theta) \right| &\leq \int \ell_t(\theta) \left| \frac{d\rho}{d\pi}(\theta) - \frac{d\rho'}{d\pi}(\theta) \right| \pi(d\theta) \\ &\leq L \underbrace{\int \left| \frac{d\rho}{d\pi}(\theta) - \frac{d\rho'}{d\pi}(\theta) \right| \pi(d\theta)}_{=2\|\rho - \rho'\|_{\text{TV}}} \end{aligned}$$

on the condition that  $0 \leq \ell_t(\theta) \leq L$  for any  $\theta$ .

## Example 2: $D_{\phi}(\rho \parallel \pi) = \chi^2(\rho \parallel \pi)$

►  $\phi(x) = x^2 - 1$  is 2-strongly convex so  $D_{\phi}$  is 2-strongly convex with respect to the  $L_2(\pi)$  norm.

► we have

$$\begin{aligned} \left| \int \ell_t(\theta) \rho(d\theta) - \int \ell_t \rho'(d\theta) \right| &\leq \int \ell_t(\theta) \left| \frac{d\rho}{d\pi}(\theta) - \frac{d\rho'}{d\pi}(\theta) \right| \pi(d\theta) \\ &\leq L \left( \int \left( \frac{d\rho}{d\pi}(\theta) - \frac{d\rho'}{d\pi}(\theta) \right)^2 \pi(d\theta) \right)^{1/2} \end{aligned}$$

on the condition that  $(\int \ell_t(\theta)^2 \pi(d\theta))^{1/2} \leq L$ .

## Constrained optimization

Constraint:  $\rho \in \mathcal{F} = \{q_\mu, \mu \in M\}$  a parametric family. Example: Gaussian distributions.

Initial objective:

$$\mu^t = \arg \min_{\mu \in M} \left\{ \sum_{s=1}^{t-1} \mathbb{E}_{\theta \sim q_\mu} [\ell_s(\theta)] + \frac{D_{\phi}(q_\mu \parallel \pi)}{\eta} \right\}.$$

Linearization gives:

$$\mu^t = \arg \min_{\mu \in M} \left\{ \sum_{s=1}^{t-1} \langle \mu, \nabla \mathbb{E}_{\theta \sim q_\mu} [\ell_s(\theta)] \rangle + \frac{D_{\phi}(q_\mu \parallel \pi)}{\eta} \right\}.$$

## Explicit update

$$F(\mu) := D_{\phi}(q_\mu \parallel \pi)$$

$$\mu_t = \nabla F^* \left( -\eta \sum_{s=1}^{t-1} \nabla_{\mu=\mu_{t-1}} \mathbb{E}_{\theta \sim q_\mu} [\ell_s(\theta)] \right).$$

Mirror descent structure: initialize  $\lambda_0 = 0$ , and update at each step:

$$\begin{cases} \lambda_t = \lambda_{t-1} - \eta \nabla_{\mu=\mu_{t-1}} \mathbb{E}_{\theta \sim q_\mu} [\ell_{t-1}(\theta)], \\ \mu_t = \nabla F^*(\lambda_t) \end{cases}$$

## Regret bound

Let  $\|\cdot\|$  be a norm on  $\mathbb{R}^d$ . If each  $\mu \mapsto \mathbb{E}_{\theta \sim q_\mu} [\ell_s(\theta)]$  is convex and  $L$ -Lipschitz with respect to  $\|\cdot\|$ , if  $\mu \mapsto D_{\phi}(q_\mu \parallel \pi)$  is  $\alpha$ -strongly convex with respect to  $\|\cdot\|$ ,

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

Example where  $D_{\phi}$  is strongly convex: Gaussian family, KL case, studied in

CHÉRIEF-ABELLATIF, B.-E., ALQUIER, P. & KHAN, M. E. (2019). A generalization bound for online variational inference. *ACML*. Conditions on the expected loss studied in

DOMKE, J. (2020). Provable smoothness guarantees for black-box variational inference. *ICML*.

The proof is based on an adaptation of the study of FTRL, see e.g.:

SHALEV-SHWARTZ, S. (2011). Online learning and online convex optimization. *Foundations and trends in Machine Learning*.