# Meta-Strategy for Learning Tuning Parameters with Guarantees

## Pierre Alquier – ABI team

**RIKEN**

**AIP**
Center for
Advanced Intelligence Project

Talk based on a joint work
with :

Dimitri Meunier

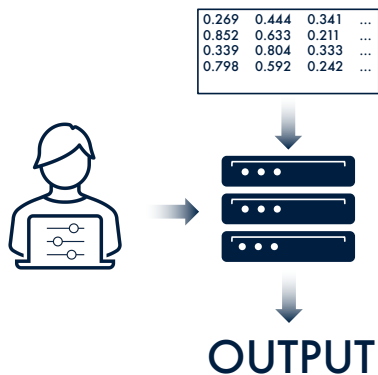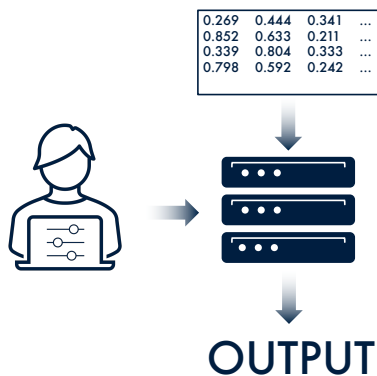2020 : ENSAE Paris and ABI team

2021 : IIT Genoa

D. Meunier, P. Alquier (2021). Meta-Strategy for Learning Tuning Parameters with Guarantees. *Preprint arXiv :2102.02504*. Submitted.

Thank you to rc3 for the drawings in this talk.
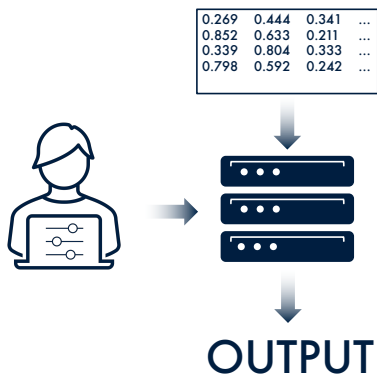
# Solving a task with an algorithm



OUTPUT

Examples :

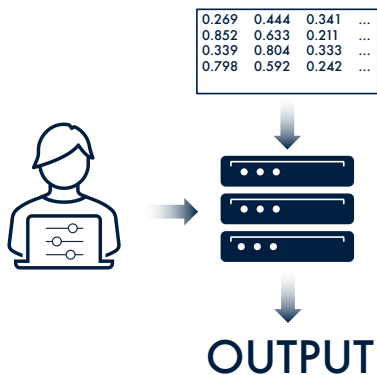| 0.269 | 0.444 | 0.341 | ... |
| 0.852 | 0.633 | 0.211 | ... |
| 0.339 | 0.804 | 0.333 | ... |
| 0.798 | 0.592 | 0.242 | ... |

**OUTPUT**

Examples :

- LASSO

$$\min_{\theta} \|y - X\theta\|^2 + \gamma \|\theta\|_1.$$

# Solving a task with an algorithm



Examples :
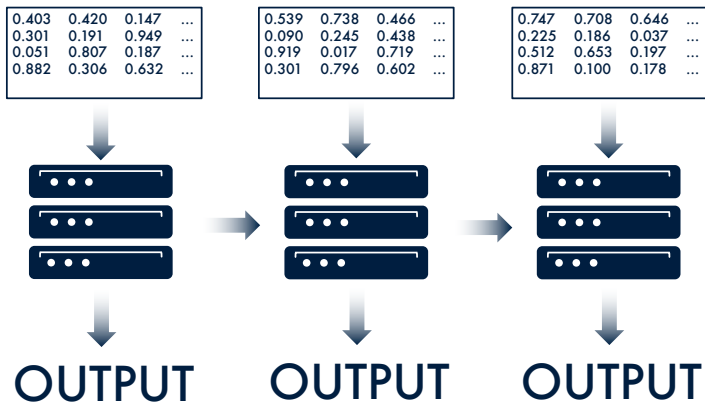
- LASSO

$$\min_\theta \|y - X\theta\|^2 + \gamma \|\theta\|_1.$$

- RIDGE

$$\min_\theta \|y - X\theta\|^2 + \alpha \|\theta\|_2^2.$$

Assume that we have an upper bound on the generalization error of a strategy when used with hyperparameter $\lambda$ on task $t$ :

$$\mathcal{L}(\mathrm{data}_t, \lambda) = \mathcal{L}_t(\lambda).$$

Assume that we have an upper bound on the generalization error of a strategy when used with hyperparameter $\lambda$ on task $t$ :

$$\mathcal{L}(\mathrm{data}_t, \lambda) = \mathcal{L}_t(\lambda).$$

Idea : use online optimization algorithm to minimize the $\mathcal{L}_t$'s...

# A meta-strategy

Assume that we have an upper bound on the generalization error of a strategy when used with hyperparameter $\lambda$ on task $t$ :

$$\mathcal{L}(\mathrm{data}_t, \lambda) = \mathcal{L}_t(\lambda).$$

Idea : use online optimization algorithm to minimize the $\mathcal{L}_t$'s...

**Online Proximal Meta-Strategy (OPMS)**

$$\lambda_{t+1} = \operatorname*{argmin}_{\lambda} \left\{ \mathcal{L}_t(\lambda) + \frac{\|\lambda - \lambda_t\|^2}{2\alpha} \right\}.$$

Sequential regression tasks $t = 1, 2, \ldots, T$

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

1. 
   - $x_{t,1}$ given

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

1. 
   - $x_{t,1}$ given
   - predict $y_{t,1}$ : $f_{\theta_{t,1}}(x_{t,1})$

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

**1**
- $x_{t,1}$ given
- predict $y_{t,1}$ : $f_{\theta_{t,1}}(x_{t,1})$
- $y_{t,1}$ is revealed

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

1. 
   - $x_{t,1}$ given
   - predict $y_{t,1} : f_{\theta_{t,1}}(x_{t,1})$
   - $y_{t,1}$ is revealed
   - update $\theta_{t,1} \rightarrow \theta_{t,2}$

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

1. - $x_{t,1}$ given
   - predict $y_{t,1}$ : $f_{\theta_{t,1}}(x_{t,1})$
   - $y_{t,1}$ is revealed
   - update $\theta_{t,1} \to \theta_{t,2}$

2. - $x_{t,2}$ given

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

1.
   - $x_{t,1}$ given
   - predict $y_{t,1}$ : $f_{\theta_{t,1}}(x_{t,1})$
   - $y_{t,1}$ is revealed
   - update $\theta_{t,1} \to \theta_{t,2}$

2.
   - $x_{t,2}$ given
   - predict $y_{t,2}$ : $f_{\theta_{t,2}}(x_{t,2})$

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

**1**
- $x_{t,1}$ given
- predict $y_{t,1} : f_{\theta_{t,1}}(x_{t,1})$
- $y_{t,1}$ is revealed
- update $\theta_{t,1} \to \theta_{t,2}$

**2**
- $x_{t,2}$ given
- predict $y_{t,2} : f_{\theta_{t,2}}(x_{t,2})$
- $y_{t,2}$ revealed

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

**1**
- $x_{t,1}$ given
- predict $y_{t,1}$ : $f_{\theta_{t,1}}(x_{t,1})$
- $y_{t,1}$ is revealed
- update $\theta_{t,1} \rightarrow \theta_{t,2}$

**2**
- $x_{t,2}$ given
- predict $y_{t,2}$ : $f_{\theta_{t,2}}(x_{t,2})$
- $y_{t,2}$ revealed
- update $\theta_{t,2} \rightarrow \theta_{t,3}$

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

**1**
- $x_{t,1}$ given
- predict $y_{t,1}$ : $f_{\theta_{t,1}}(x_{t,1})$
- $y_{t,1}$ is revealed
- update $\theta_{t,1} \to \theta_{t,2}$

**2**
- $x_{t,2}$ given
- predict $y_{t,2}$ : $f_{\theta_{t,2}}(x_{t,2})$
- $y_{t,2}$ revealed
- update $\theta_{t,2} \to \theta_{t,3}$

**3**
- $x_{t,3}$ given

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

1.
   - $x_{t,1}$ given
   - predict $y_{t,1} : f_{\theta_{t,1}}(x_{t,1})$
   - $y_{t,1}$ is revealed
   - update $\theta_{t,1} \to \theta_{t,2}$

2.
   - $x_{t,2}$ given
   - predict $y_{t,2} : f_{\theta_{t,2}}(x_{t,2})$
   - $y_{t,2}$ revealed
   - update $\theta_{t,2} \to \theta_{t,3}$

3.
   - $x_{t,3}$ given
   - predict $y_{t,3} : f_{\theta_{t,3}}(x_{t,3})$

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

**1**
- $x_{t,1}$ given
- predict $y_{t,1} : f_{\theta_{t,1}}(x_{t,1})$
- $y_{t,1}$ is revealed
- update $\theta_{t,1} \to \theta_{t,2}$

**2**
- $x_{t,2}$ given
- predict $y_{t,2} : f_{\theta_{t,2}}(x_{t,2})$
- $y_{t,2}$ revealed
- update $\theta_{t,2} \to \theta_{t,3}$

**3**
- $x_{t,3}$ given
- predict $y_{t,3} : f_{\theta_{t,3}}(x_{t,3})$
- $y_{t,3}$ revealed

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

**1**
- $x_{t,1}$ given
- predict $y_{t,1} : f_{\theta_{t,1}}(x_{t,1})$
- $y_{t,1}$ is revealed
- update $\theta_{t,1} \to \theta_{t,2}$

**2**
- $x_{t,2}$ given
- predict $y_{t,2} : f_{\theta_{t,2}}(x_{t,2})$
- $y_{t,2}$ revealed
- update $\theta_{t,2} \to \theta_{t,3}$

**3**
- $x_{t,3}$ given
- predict $y_{t,3} : f_{\theta_{t,3}}(x_{t,3})$
- $y_{t,3}$ revealed
- update $\theta_{t,3} \to \theta_{t,4}$

**4** $\ldots$

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

**1**
- $x_{t,1}$ given
- predict $y_{t,1}$ : $f_{\theta_{t,1}}(x_{t,1})$
- $y_{t,1}$ is revealed
- update $\theta_{t,1} \to \theta_{t,2}$

Loss function $\ell$.

**2**
- $x_{t,2}$ given
- predict $y_{t,2}$ : $f_{\theta_{t,2}}(x_{t,2})$
- $y_{t,2}$ revealed
- update $\theta_{t,2} \to \theta_{t,3}$

**3**
- $x_{t,3}$ given
- predict $y_{t,3}$ : $f_{\theta_{t,3}}(x_{t,3})$
- $y_{t,3}$ revealed
- update $\theta_{t,3} \to \theta_{t,4}$

**4** $\ldots$

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

1. - $x_{t,1}$ given
   - predict $y_{t,1} : f_{\theta_{t,1}}(x_{t,1})$
   - $y_{t,1}$ is revealed
   - update $\theta_{t,1} \to \theta_{t,2}$

2. - $x_{t,2}$ given
   - predict $y_{t,2} : f_{\theta_{t,2}}(x_{t,2})$
   - $y_{t,2}$ revealed
   - update $\theta_{t,2} \to \theta_{t,3}$

3. - $x_{t,3}$ given
   - predict $y_{t,3} : f_{\theta_{t,3}}(x_{t,3})$
   - $y_{t,3}$ revealed
   - update $\theta_{t,3} \to \theta_{t,4}$

4. $\ldots$

Loss function $\ell$.
For short,

$$\ell_{t,i}(\theta) = \ell(y_{t,i}, f_\theta(x_{t,i})).$$

# Example : tasks are online prediction

## Sequential regression tasks $t = 1, 2, \ldots, T$

**1**
- $x_{t,1}$ given
- predict $y_{t,1} : f_{\theta_{t,1}}(x_{t,1})$
- $y_{t,1}$ is revealed
- update $\theta_{t,1} \to \theta_{t,2}$

**2**
- $x_{t,2}$ given
- predict $y_{t,2} : f_{\theta_{t,2}}(x_{t,2})$
- $y_{t,2}$ revealed
- update $\theta_{t,2} \to \theta_{t,3}$

**3**
- $x_{t,3}$ given
- predict $y_{t,3} : f_{\theta_{t,3}}(x_{t,3})$
- $y_{t,3}$ revealed
- update $\theta_{t,3} \to \theta_{t,4}$

**4** $\ldots$

Loss function $\ell$.
For short,

$$\ell_{t,i}(\theta) = \ell(y_{t,i}, f_\theta(x_{t,i})).$$

Objective :

$$\sum_{i=1}^{n} \ell_{t,i}(\theta_{t,i})$$

as small as possible.

# Example : online gradient algorithm (OGA)

$$\theta_{t,i+1} = \theta_{t,i} - \eta \nabla \ell_{t,i}(\theta_{t,i})$$

# Example : online gradient algorithm (OGA)

$$\theta_{t,i+1} = \theta_{t,i} - \eta \nabla \ell_{t,i}(\theta_{t,i})$$
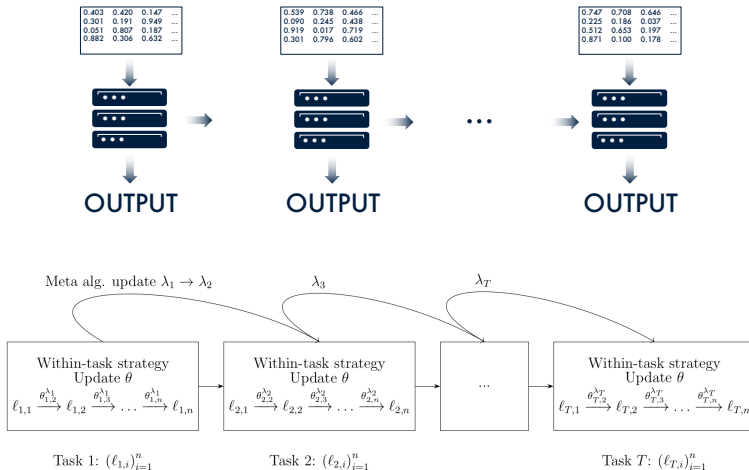
### Regret bound for OGA

If each $\ell_{t,i}$ is convex and $L$-Lipschitz,

$$\sum_{i=1}^{n} \ell_{t,i}(\theta_{t,i}) \leq \underbrace{\inf_{\|\theta\| \leq B} \left\{ \sum_{i=1}^{n} \ell_{t,i}(\theta) + \frac{\eta n L^2}{2} + \frac{\|\theta - \theta_{1,t}\|^2}{2\eta} \right\}}_{\mathcal{L}_t(\eta,\theta_{t,1})=\mathcal{L}_t(\lambda)}.$$

# Example : online gradient algorithm (OGA)

$$\theta_{t,i+1} = \theta_{t,i} - \eta \nabla \ell_{t,i}(\theta_{t,i})$$

## Regret bound for OGA

If each $\ell_{t,i}$ is convex and $L$-Lipschitz,

$$\sum_{i=1}^{n} \ell_{t,i}(\theta_{t,i}) \leq \inf_{\|\theta\| \leq B} \underbrace{\left\{ \sum_{i=1}^{n} \ell_{t,i}(\theta) + \frac{\eta n L^2}{2} + \frac{\|\theta - \theta_{1,t}\|^2}{2\eta} \right\}}_{\mathcal{L}_t(\eta, \theta_{t,1}) = \mathcal{L}_t(\lambda)}.$$

Take $\eta \sim 1/\sqrt{n}$ to get :

$$\sum_{i=1}^{n} \ell_{t,i}(\theta_{t,i}) \leq \inf_{\|\theta\| \leq B} \quad \sum_{i=1}^{n} \ell_{t,i}(\theta) + \mathcal{O}(\sqrt{n}).$$

# Example : online gradient algorithm (OGA)

$$\theta_{t,i+1} = \theta_{t,i} - \eta \nabla \ell_{t,i}(\theta_{t,i})$$

## Regret bound for OGA

If each $\ell_{t,i}$ is convex and $L$-Lipschitz,

$$\sum_{i=1}^{n} \ell_{t,i}(\theta_{t,i}) \leq \inf_{\|\theta\| \leq B} \underbrace{\left\{ \sum_{i=1}^{n} \ell_{t,i}(\theta) + \frac{\eta n L^2}{2} + \frac{\|\theta - \theta_{1,t}\|^2}{2\eta} \right\}}_{\mathcal{L}_t(\eta, \theta_{t,1}) = \mathcal{L}_t(\lambda)}.$$

Take $\eta \sim 1/\sqrt{n}$ to get :

$$\frac{1}{n} \sum_{i=1}^{n} \ell_{t,i}(\theta_{t,i}) \leq \inf_{\|\theta\| \leq B} \frac{1}{n} \sum_{i=1}^{n} \ell_{t,i}(\theta) + \mathcal{O}\left( \frac{1}{\sqrt{n}} \right).$$

# Meta-learning for OGA

# Meta-learning for OGA

## Online Proximal Meta-Strategy (OPMS)

$$\lambda_{t+1} = \operatorname*{argmin}_{\lambda} \left\{ \mathcal{L}_t(\lambda) + \frac{\|\lambda - \lambda_t\|^2}{2\alpha} \right\}.$$

# Meta-learning for OGA

**Online Proximal Meta-Strategy (OPMS)**

$$\lambda_{t+1} = \operatorname*{argmin}_{\lambda} \left\{ \mathcal{L}_t(\lambda) + \frac{\|\lambda - \lambda_t\|^2}{2\alpha} \right\}.$$

$(\eta_{t+1}, \theta_{t+1,1})$

$$= \operatorname*{argmin}_{\eta, \vartheta} \min_{\|\theta\| \leq B} \left\{ \sum_{i=1}^{n} \ell_{t,i}(\theta) + \frac{\eta n L^2}{2} + \frac{\|\theta - \vartheta\|^2}{2\eta} \right.$$

$$\left. + \frac{\|\vartheta - \theta_{t,1}\|^2 + (\eta - \eta_t)^2}{2\alpha} \right\}.$$

Question : what do we win ?

# Standard : learning in isolation

$$\frac{1}{n} \sum_{i=1}^{n} \ell_{1,i}(\theta_{1,i}) \leq \inf_{\|\theta\| \leq B} \frac{1}{n} \sum_{i=1}^{n} \ell_{1,i}(\theta) + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$$
$$+$$
$$\vdots$$
$$+$$
$$\frac{1}{n} \sum_{i=1}^{n} \ell_{T,i}(\theta_{T,i}) \leq \inf_{\|\theta\| \leq B} \frac{1}{n} \sum_{i=1}^{n} \ell_{T,i}(\theta) + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$$
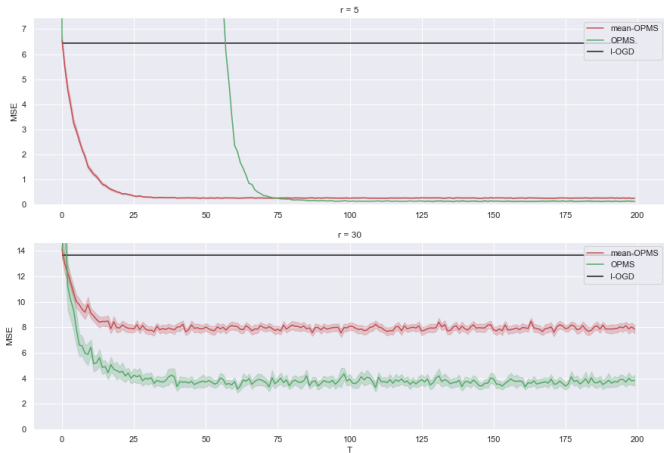
# Standard : learning in isolation

$$\frac{1}{n} \sum_{i=1}^{n} \ell_{1,i}(\theta_{1,i}) \leq \inf_{\|\theta\| \leq B} \frac{1}{n} \sum_{i=1}^{n} \ell_{1,i}(\theta) + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$$

$$+$$

$$\vdots$$

$$+$$

$$\frac{1}{n} \sum_{i=1}^{n} \ell_{T,i}(\theta_{T,i}) \leq \inf_{\|\theta\| \leq B} \frac{1}{n} \sum_{i=1}^{n} \ell_{T,i}(\theta) + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$$
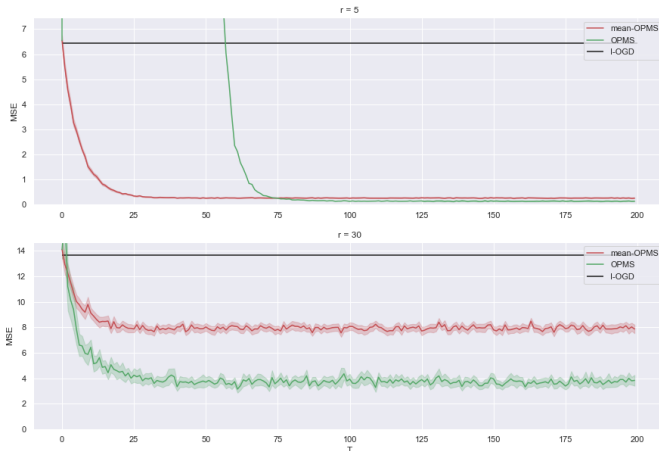
$$\frac{1}{nT} \sum_{t=1}^{T} \sum_{i=1}^{n} \ell_{t,i}(\theta_{1,i}) \leq \inf_{\substack{\|\theta_1\| \leq B \\ \cdots \\ \|\theta_T\| \leq B}} \frac{1}{nT} \sum_{t=1}^{T} \sum_{i=1}^{n} \ell_{t,i}(\theta_t) + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right).$$

**Theorem**

$$\frac{1}{nT} \sum_{t=1}^{T} \sum_{i=1}^{n} \ell_{t,i}(\theta_{1,i}) \leq \inf_{\substack{\|\theta_{\mathbf{1}}\| \leq B \\ \cdots \\ \|\theta_T\| \leq B}} \left\{ \frac{1}{nT} \sum_{t=1}^{T} \sum_{i=1}^{n} \ell_{t,i}(\theta_t) \right.$$

$$\left. + \mathcal{O}\left( \frac{\sqrt{\frac{1}{T} \sum_{t=1}^{T} \left(\theta_t - \bar{\theta}\right)^2}}{\sqrt{n}} + \frac{1}{n} + \frac{n}{\sqrt{T}} \right) \right\}.$$

# Simulated examples

# Simulated examples

# But....

## Approximate Bayesian Inference Team



Mohammad Emtiyaz Khan
(Ph.D.)

Title

Team Leader

# ximate Bayesian Inferen

**Title**

Team Lea
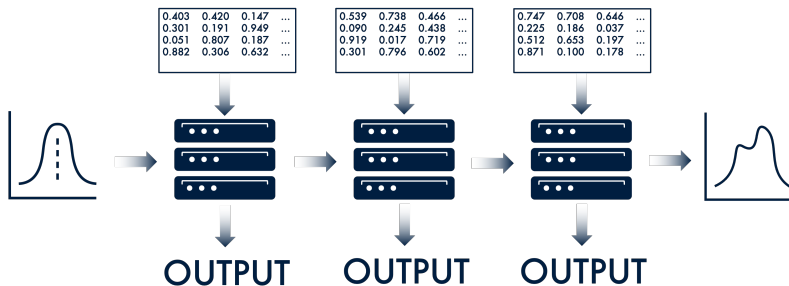
Bayesian

# Learning the prior

# Online variational inference

B.-E. Chérief-Abdellatif, P. Alquier, M. E. Khan (2019). *A generalization bound for online variational inference*. ACML.

Online variational inference :

$$
\mu_{t,i} = \operatorname*{argmin}_{\mu \in M} \left\{ \mu^T \sum_{j=1}^{i-1} \nabla_{\mu_{t,j}} \mathbb{E}_{\theta \sim q_{\mu_{t,j}}} \left[ \ell_{t,j}(\theta) \right] + \frac{\mathcal{K}(q_\mu, \pi)}{\eta} \right\}.
$$

# Online variational inference

B.-E. Chérief-Abdellatif, P. Alquier, M. E. Khan (2019). *A generalization bound for online variational inference*. ACML.

Online variational inference :

$$\mu_{t,i} = \underset{\mu \in M}{\operatorname{argmin}} \left\{ \mu^T \sum_{j=1}^{i-1} \nabla_{\mu_{t,j}} \mathbb{E}_{\theta \sim q_{\mu_{t,j}}} \left[ \ell_{t,j}(\theta) \right] + \frac{\mathcal{K}(q_\mu, \pi)}{\eta} \right\}.$$

Regret bound :

$$\sum_{i=1}^{n} \mathbb{E}_{\theta \sim q_{\mu_{t,i}}} [\ell_{t,i}(\theta)] \leq \inf_{\mu \in \mathcal{M}} \left\{ \mathbb{E}_{\theta \sim q_\mu} \left[ \sum_{i=1}^{n} \ell_{t,i}(\theta) \right] + \frac{\eta 4 L^2 n}{\alpha} + \frac{\mathcal{K}(q_\mu, \pi)}{\eta} \right\}.$$

# Online variational inference

B.-E. Chérief-Abdellatif, P. Alquier, M. E. Khan (2019). *A generalization bound for online variational inference*. ACML.

Online variational inference :

$$\mu_{t,i} = \operatorname*{argmin}_{\mu \in M} \left\{ \mu^T \sum_{j=1}^{i-1} \nabla_{\mu_{t,j}} \mathbb{E}_{\theta \sim q_{\mu_{t,j}}} [\ell_{t,j}(\theta)] + \frac{\mathcal{K}(q_\mu, \pi)}{\eta} \right\}.$$

Regret bound :

$$\sum_{i=1}^n \mathbb{E}_{\theta \sim q_{\mu_{t,i}}} [\ell_{t,i}(\theta)] \leq \inf_{\mu \in \mathcal{M}} \left\{ \mathbb{E}_{\theta \sim q_\mu} \left[ \sum_{i=1}^n \ell_{t,i}(\theta) \right] + \frac{\eta 4 L^2 n}{\alpha} + \frac{\mathcal{K}(q_\mu, \pi)}{\eta} \right\}.$$

If $q_\mu = \mathcal{N}(\mu, I)$ and $\pi = \mathcal{N}(m, I)$, $\mathcal{K}(q_\mu, \pi) = \frac{\|\mu - m\|^2}{2}$.

どうも ありがとう ございました!