

# Quantized Variational Inference

## Optimal Quantization for ELBO maximization

Amir DIB

ENS Paris-Saclay, Centre Borelli, Université Paris-Saclay, CNRS.

### MOTIVATION

**Quantized Variational Inference** is a new algorithm for Evidence Lower Bound maximization. Optimal **Voronoi Tessellation** produces **variance free** gradients for ELBO optimization at the cost of introducing asymptotically decaying bias. Using the **Quantized Variational Inference framework** leads to **fast convergence** for both score function and the reparametrized gradient estimator at a comparable computational cost.

### OVERVIEW

Let  $y$  be the data,  $z$  a latent variables and  $p(y, z)$  the model. The goal of the Bayesian statistician is to find the best latent variable that fits the data, hence the likelihood  $p(z|y)$ . We can approximate the true posterior by maximizing the **Evidence Lower Bound**  $\mathcal{L}(\lambda)$  though the variational distribution  $q_\lambda$  thanks to the following decomposition

$$\log p(y) = \underbrace{\mathbb{E}_{z \sim q_\lambda} \left[ \log \frac{p(z, y)}{q_\lambda(z)} \right]}_{\mathcal{L}(\lambda)} + \text{KL} (q_\lambda(z) \| p(z|y)).$$

Given a sample  $(X_1^\lambda, \dots, X_N^\lambda)$  of size  $N$ , typical MCVI procedure consists of a Gradient descent at each step  $k$

$$\lambda_{k+1} = \lambda_k - \alpha_k \underbrace{\frac{1}{N} \sum_{i=1}^N \nabla_\lambda f(X_i^{\lambda_k})}_{\hat{g}_{MC}^N}.$$

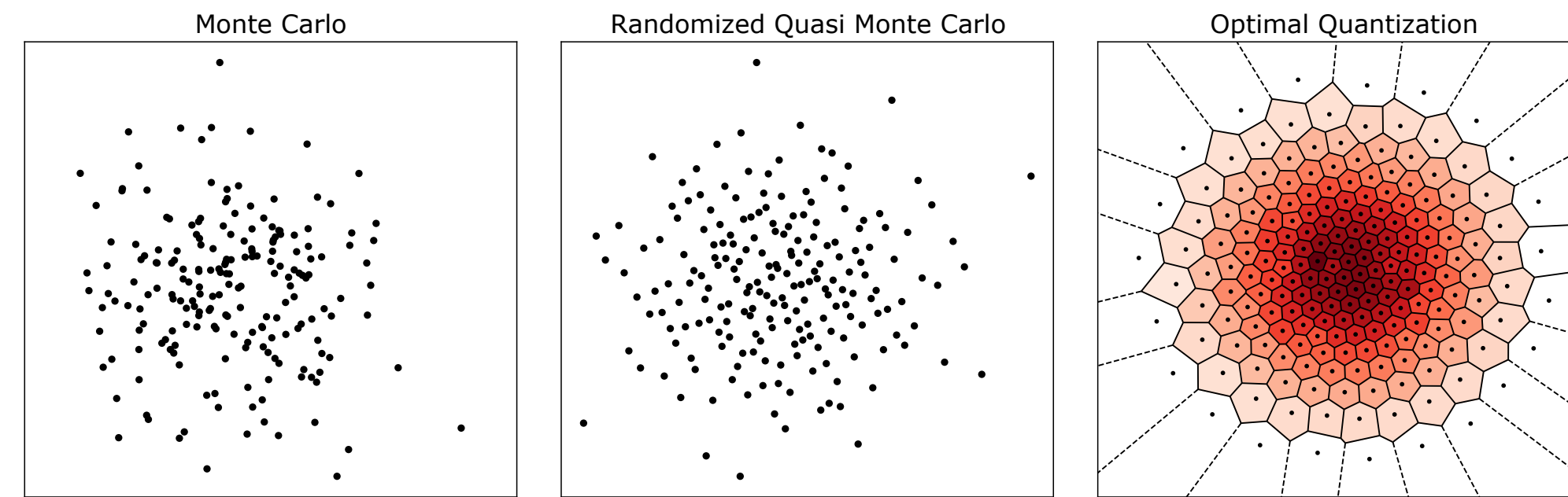
The speed of convergence crucially depends on the expected norm of the gradient and thus on the gradient variance thanks to the bias-variance decomposition

$$\mathbb{E} \|g\|_{\ell_2}^2 = \text{tr} \mathbb{V}g + \|\mathbb{E}g\|_{\ell_2}^2.$$

Our approach consists of considering alternative sampling instead of the traditional Monte Carlo. Precisely, we consider the **Optimal Quantizer** or **Voronoi Tessellation** of  $q_\lambda$  at level  $N$ ,  $X^{\Gamma_N, \lambda}$ . We thus obtain a **deterministic gradient** allowing for larger step in the optimization procedure by using the following gradient descent scheme

$$\lambda_{k+1} = \lambda_k - \alpha_k \underbrace{\nabla_\lambda \sum_{i=1}^N \omega_i^k f(X_i^{\Gamma_N, \lambda_k})}_{\hat{g}_{OQ}^N}.$$

### OPTIMAL QUANTIZER



Different type of sampling for a bi-variate standard gaussian distribution.

Take  $\Gamma_N = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , the  $L^p_{\mathbb{R}^d}$  quantiser is defined as the probability measure on the convex subset of probability measure on  $\Gamma_N$  that minimizes the Wasserstein distance

$$\inf_{\Gamma \in \mathbb{R}^d} \int_{\mathbb{R}^d} \min_{1 \leq k \leq K} |\xi - x_k|^p \mu(d\xi).$$

Any expectation can be computed with the cubature formula

$$\mathbb{E}f(X^{\Gamma_N}) = \sum_{i=1}^N \mathbb{P}(X^{\Gamma_N} = x_i) f(x_i).$$

Denoting  $\hat{\mathcal{L}}_{OQ}^N(\lambda)$  the ELBO obtained by sampling with the **OQ** and  $\mathcal{L}(\lambda)$  the true ELBO, we can obtain a bound on the produced bias

$$\left| \mathcal{L}(\lambda) - \hat{\mathcal{L}}_{OQ}^N(\lambda) \right| \leq C \|X^\lambda - X^{\Gamma_N, \lambda}\|_2.$$

### Quantized Variational Inference

**Input:**  $y, p(x, z), q_{\lambda_0}$ .

**Result:** Optimal Quantized VI parameters  $\lambda_q^*$ .

**while not converged do**

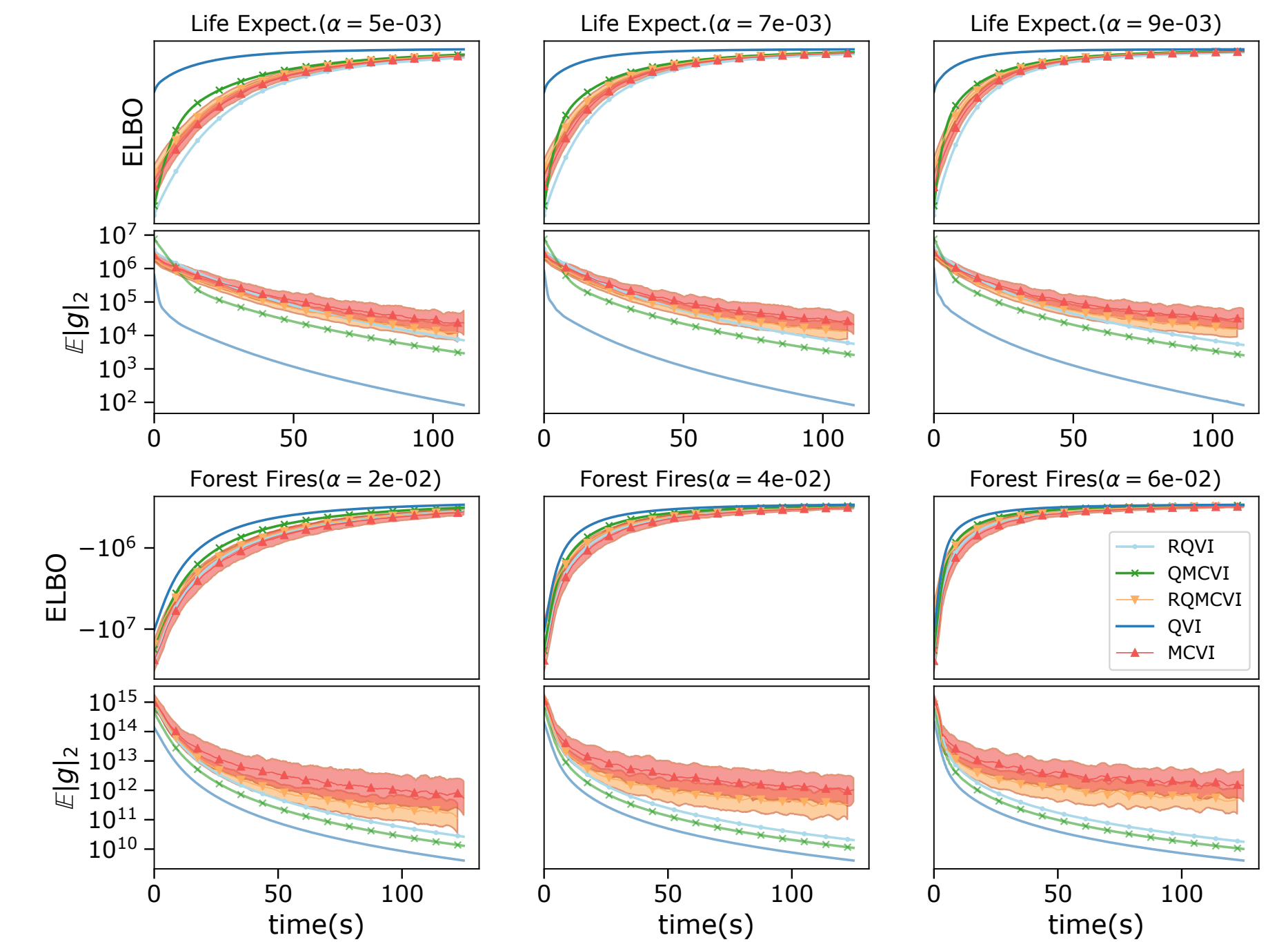
    Get  $(X_1^{\Gamma_N, \lambda_k}, \dots, X_N^{\Gamma_N, \lambda_k}) \sim q_{\lambda_k}, (w_1^k, \dots, w_N^k)$ ;

    Compute  $\hat{g}_{OQ}^N(\lambda_k) = \nabla_\lambda \sum_{i=1}^N w_i^k H(X_i^{\Gamma_N, \lambda_k})$ ;

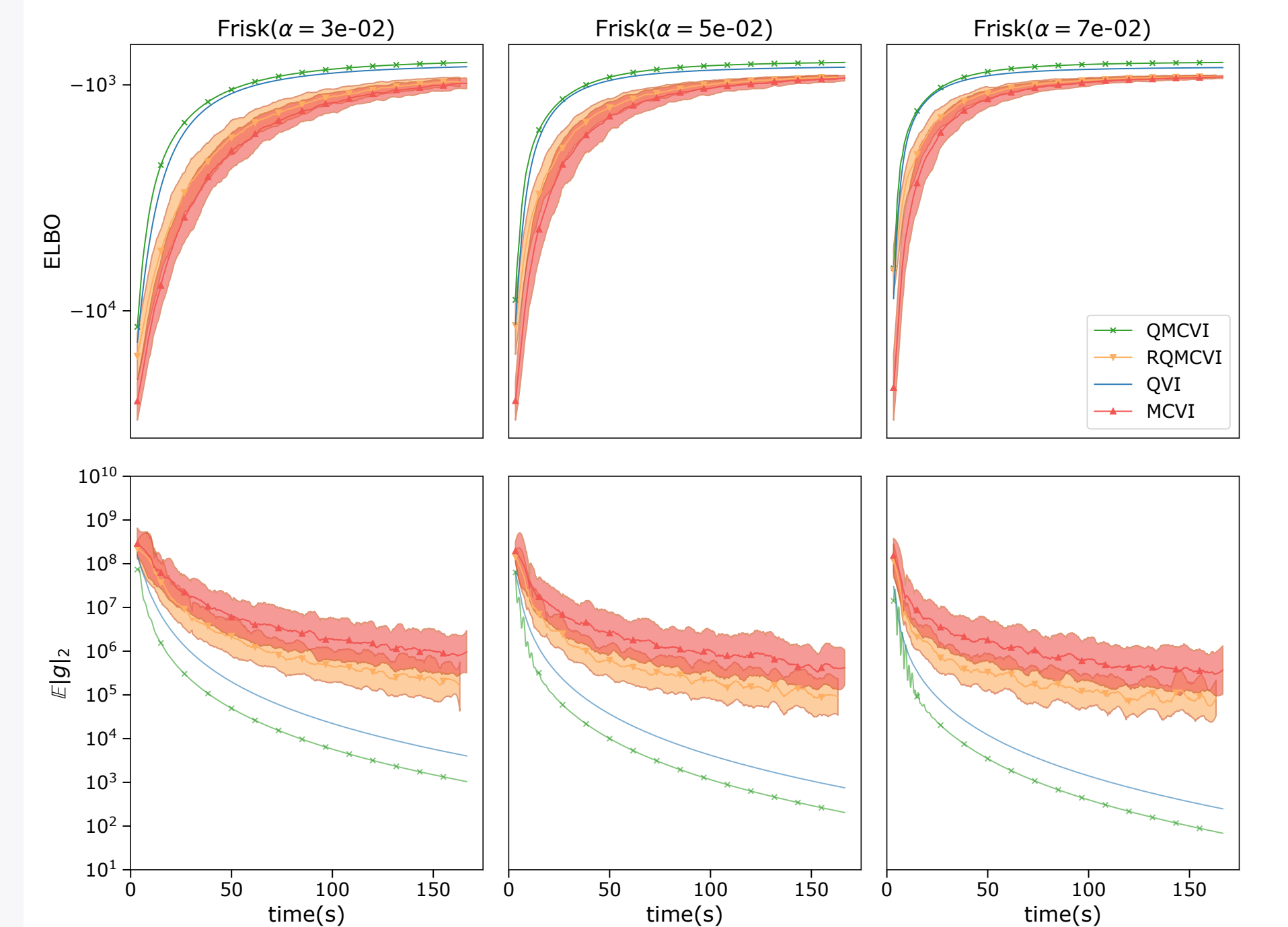
$\lambda_{k+1} = \lambda_k - \alpha_k \hat{g}_{OQ}^N(\lambda_k)$ ;

**end**

### EXPERIMENTS



Bayesian Linear Regression.



Hierarchical Poisson Model.

#### References

- Gilles Pagès (2018). "Numerical Probability: An Introduction with Application to Finance". Springer International Publishing
- Léon Bottou (2018), Frank E. Curtis, and Jorge Nocedal. "Optimization Methods for Large-Scale Machine Learning". In: SIAM Review 60.2, pp. 223–311.

