Adaptive Importance Sampling meets Mirror Descent: a Bias-variance tradeoff

> Anna Korba¹ François Portier² ¹ENSAE, CREST, ²ENSAI, CREST

> > JSTAR Rennes 2022

Outline

Problem

Regularized Importance Sampling

Safe and Regularized Adaptive Importance Sampling

Uniform convergence of the scheme

Adaptive choice of the regularization parameter

Experiments

Main motivation = Bayesian inference:

Main motivation = Bayesian inference:

• Let $\mathcal{D} = (x_i, y_i)_{i=1}^m$ a labelled dataset of i.i.d. points.

Main motivation = Bayesian inference:

• Let $\mathcal{D} = (x_i, y_i)_{i=1}^m$ a labelled dataset of i.i.d. points.

▶ Assume an underlying model parametrized by $z \in \mathbb{R}^d$, e.g.

$$y = f(x, z) + \epsilon$$

(p(y|x,z) gaussian)

Main motivation = Bayesian inference:

- Let $\mathcal{D} = (x_i, y_i)_{i=1}^m$ a labelled dataset of i.i.d. points.
- ▶ Assume an underlying model parametrized by $z \in \mathbb{R}^d$, e.g.

$$y = f(x, z) + \epsilon$$

(p(y|x, z) gaussian)

Compute the likelihood:

$$p(\mathcal{D}|z) = \prod_{i=1}^{m} p(y_i|x_i, z)$$

Main motivation = Bayesian inference:

- Let $\mathcal{D} = (x_i, y_i)_{i=1}^m$ a labelled dataset of i.i.d. points.
- ▶ Assume an underlying model parametrized by $z \in \mathbb{R}^d$, e.g.

$$y = f(x, z) + \epsilon$$

(p(y|x, z) gaussian)

Compute the likelihood:

$$p(\mathcal{D}|z) = \prod_{i=1}^{m} p(y_i|x_i, z)$$

• Assume a prior distribution on the parameter $z \sim p$.

Bayes' rule :
$$f(z) := p(z|\mathcal{D}) = \frac{p(\mathcal{D}|z)p(z)}{C}, C = \int_{\mathbb{R}^d} p(\mathcal{D}|z)p(z)dz.$$

f is known up to a constant since C is intractable.

How to sample from *f* then? e.g. to compute the "Bayesian model average":

$$p(y|x, \mathcal{D}) = \int_{\mathbb{R}^d} p(y|x, z) df(z)$$

Bayes' rule :
$$f(z) := p(z|\mathcal{D}) = \frac{p(\mathcal{D}|z)p(z)}{C}, C = \int_{\mathbb{R}^d} p(\mathcal{D}|z)p(z)dz.$$

f is known up to a constant since C is intractable.

How to sample from *f* then? e.g. to compute the "Bayesian model average":

$$p(y|x, \mathcal{D}) = \int_{\mathbb{R}^d} p(y|x, z) df(z)$$

- 1. **MCMC methods** (Markov Chain Monte Carlo): generate a markov chain $(X_t)_{t\geq 0}$ whose law q_t converges to f as $t \to \infty$
- 2. Variational Inference

$$\tilde{f} = \operatorname*{argmin}_{q \in \mathcal{Q}} \operatorname{KL}(q|f)$$

where $\ensuremath{\mathcal{Q}}$ is a parametric family of probability distributions

3. **Importance sampling** : sample from a simple proposal distribution *q* dominating *f* and reweight

Naive Importance Sampling

Let X a random variable with distribution q dominating f. The basic idea of IS is to re-weight g(X) by **the importance weight** W(X) = f(X)/q(X).

Naive Importance Sampling

Let X a random variable with distribution q dominating f. The basic idea of IS is to re-weight g(X) by the importance weight W(X) = f(X)/q(X).

Since $\mathbb{E}[W(X)g(X)] = \int gf$ and using i.i.d. samples $X_1, \ldots, X_n \sim q$, one can build an (unbiased) IS estimator of $\int gf$ as

$$\int gf \approx \frac{1}{n} \sum_{k=1}^n \frac{f(X_k)}{q(X_k)} g(X_k) = \frac{1}{n} \sum_{k=1}^n W(X_k) g(X_k).$$

Naive Importance Sampling

Let X a random variable with distribution q dominating f. The basic idea of IS is to re-weight g(X) by **the importance weight** W(X) = f(X)/q(X).

Since $\mathbb{E}[W(X)g(X)] = \int gf$ and using i.i.d. samples $X_1, \ldots, X_n \sim q$, one can build an (unbiased) IS estimator of $\int gf$ as

$$\int gf \approx \frac{1}{n} \sum_{k=1}^n \frac{f(X_k)}{q(X_k)} g(X_k) = \frac{1}{n} \sum_{k=1}^n W(X_k) g(X_k).$$

Remark: if *f* is known up to a normalization constant, use normalized weights $\sum_{k=1}^{n} W(X_k)g(X_k) / \sum_{k=1}^{n} W(X_k)$.

Naive importance sampling

Problem: if *q* is far from the target *f*, the importance weights may have a large variance (hence the IS estimator as well) !



The accuracy heavily depends on the choice of q

Contributions of the paper

Approach: Adaptive Importance Sampling (AIS)

Idea: use a sequence of proposals $(q_k)_{k\geq 0}$, learning from the past actions and data, to approximate *f*.

Contributions of the paper

Approach: Adaptive Importance Sampling (AIS)

Idea: use a sequence of proposals $(q_k)_{k\geq 0}$, learning from the past actions and data, to approximate *f*.

We propose a new non parametric AIS method, that

- (i) introduces a new regularization strategy which raises adaptively the importance sampling weights to a certain power ranging from 0 to 1
- (ii) uses a mixture between a kernel density estimate of the target and a safe reference density as proposal.

Outline

Problem

Regularized Importance Sampling

Safe and Regularized Adaptive Importance Sampling

Uniform convergence of the scheme

Adaptive choice of the regularization parameter

Experiments

Idea: use regularized weights of the form $W(X)^{\eta}$, $\eta \in (0, 1)$.

Idea: use regularized weights of the form $W(X)^{\eta}$, $\eta \in (0, 1)$.

Lemma: Suppose that *q* dominates *f* and define W(X) = f(X)/q(X) with *X* having density *q*. For all $\eta \in (0, 1]$: $\mathbb{E}[W(X)^{\eta}] \leq 1$ and $\operatorname{Var}[W(X)^{\eta}] \leq \operatorname{Var}[W(X)]$.

Idea: use regularized weights of the form $W(X)^{\eta}$, $\eta \in (0, 1)$.

Lemma: Suppose that *q* dominates *f* and define W(X) = f(X)/q(X) with *X* having density *q*. For all $\eta \in (0, 1]$: $\mathbb{E}[W(X)^{\eta}] \leq 1$ and $Var[W(X)^{\eta}] \leq Var[W(X)]$.

Proof:

- Jensen's inequality: E[W(X)^η] ≤ E[W(X)]^η = 1 since f ≪ q.
- we have, since $|w^{\eta} 1| \le |w 1|$ for all $w \ge 0$:

$$\begin{aligned} \mathsf{Var}[W(X)^{\eta}] &\leq \mathsf{Var}[W(X)^{\eta}] + (\mathbb{E}[W(X)^{\eta}] - 1)^2 \\ &= \mathbb{E}[(W(X)^{\eta} - 1)^2] \leq \mathbb{E}[(W(X) - 1)^2]. \end{aligned}$$

Remarks:

• choosing η enables to balance bias and variance !

•
$$\mathbb{E}[W(X)^{\eta}g(X)] = \int f^{\eta}q^{1-\eta}g$$

Hence, regularized IS "moves" from the initial density q ($\eta = 0$) to the target density $f^{\eta}q^{1-\eta}$ (=*f* if $\eta = 1$).

Remarks:

• choosing η enables to balance bias and variance !

•
$$\mathbb{E}[W(X)^{\eta}g(X)] = \int f^{\eta}q^{1-\eta}g$$

Hence, regularized IS "moves" from the initial density q ($\eta = 0$) to the target density $f^{\eta}q^{1-\eta}$ (=*f* if $\eta = 1$).

Additional Remarks:

different from simulated annealing (sequence of tempered posteriors) : *f* → *f^η* ⇒ would yield IS weights of the form *f^η/q* instead of (*f/q*)^η

• can be seen as (entropic) mirror descent with step-size η_k :

$$q_{k+1} \propto q_k^{1-\eta_k} f^{\eta_k}$$

Consider the sampling/variational inference objective:

$$q^{\star} = \operatorname*{arginf}_{q \in \mathcal{Q}} \ \operatorname{\mathsf{KL}}(q|f), \quad \operatorname{\mathsf{KL}}(q|f) = \int_{\mathbb{R}^d} \log\Bigl(rac{q}{f}(x)\Bigr) dq(x)$$

Consider the sampling/variational inference objective:

$$q^{\star} = \operatorname*{arginf}_{q \in \mathcal{Q}} \ \operatorname{\mathsf{KL}}(q|f), \ \ \operatorname{\mathsf{KL}}(q|f) = \int_{\mathbb{R}^d} \log\Bigl(rac{q}{f}(x)\Bigr) dq(x)$$

Entropic mirror descent applied to this objective can be written at each time $k \ge 0$:

$$q_{k+1}^{*} = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \quad \eta_{k} \underbrace{\int_{\mathbb{R}^{d}} \log\left(\frac{q_{k}^{*}}{f}(x)\right) d(q - q_{k}^{*})(x)}_{\langle \nabla \operatorname{KL}(q_{k}^{*}|f), q - q_{k}^{*} \rangle} + \operatorname{KL}(q|q_{k}^{*}).$$
(1)

Consider the sampling/variational inference objective:

$$q^{\star} = \operatorname*{arginf}_{q \in \mathcal{Q}} \ \operatorname{\mathsf{KL}}(q|f), \ \ \operatorname{\mathsf{KL}}(q|f) = \int_{\mathbb{R}^d} \log\Bigl(rac{q}{f}(x)\Bigr) dq(x)$$

Entropic mirror descent applied to this objective can be written at each time $k \ge 0$:

$$q_{k+1}^{*} = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \quad \eta_{k} \underbrace{\int_{\mathbb{R}^{d}} \log\left(\frac{q_{k}^{*}}{f}(x)\right) d(q - q_{k}^{*})(x)}_{\langle \nabla \operatorname{KL}(q_{k}^{*}|f), q - q_{k}^{*} \rangle} + \operatorname{KL}(q|q_{k}^{*}).$$
(1)

Differentiate (1) w.r.t. q yields:

$$\eta_k \log\left(\frac{q_k^*}{f}\right) + \log\left(\frac{q_{k+1}^*}{q_k^*}\right) = 0$$

Consider the sampling/variational inference objective:

$$q^{\star} = \operatorname*{arginf}_{q \in \mathcal{Q}} \ \operatorname{\mathsf{KL}}(q|f), \ \ \operatorname{\mathsf{KL}}(q|f) = \int_{\mathbb{R}^d} \log\Bigl(rac{q}{f}(x)\Bigr) dq(x)$$

Entropic mirror descent applied to this objective can be written at each time $k \ge 0$:

$$q_{k+1}^{*} = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \quad \eta_{k} \underbrace{\int_{\mathbb{R}^{d}} \log\left(\frac{q_{k}^{*}}{f}(x)\right) d(q - q_{k}^{*})(x)}_{\langle \nabla \operatorname{KL}(q_{k}^{*}|f), q - q_{k}^{*} \rangle} + \operatorname{KL}(q|q_{k}^{*}).$$
(1)

Differentiate (1) w.r.t. q yields:

$$\eta_k \log\left(rac{q_k^*}{f}
ight) + \log\left(rac{q_{k+1}^*}{q_k^*}
ight) = 0$$

Assuming $Q = \mathcal{P}(\mathbb{R}^d)$, (1) is minimized for:

$$q_{k+1}^*(x) = \frac{f(x)^{\eta_k} q_k^*(x)^{1-\eta_k}}{\int_{\mathbb{R}^d} f(x')^{\eta_k} q_k^*(x')^{1-\eta_k} dx'} \propto f^{\eta_k}(x) q_k^{*(1-\eta_k)}(x).$$

Fast convergence of Entropic Mirror Descent (EMD)

The previous scheme is attractive thanks to its fast convergence.

Lemma

Let $(\eta_k)_{k\geq 1}$ valued in (0, 1] and $(q_k^*)_{k\geq 1}$ be defined by EMD starting from an initial probability density function q_1 . Then, for all $n \in \mathbb{N}^*$,

$$\int_{\mathbb{R}^d} |f(x) - q_{n+1}^*(x)| dx \leq \sqrt{2 \operatorname{KL}(f|q_1)} \prod_{k=1}^n (1 - \eta_k)^{1/2},$$

• $\eta_k = c/k$ with 0 < c < 1 yields the rate $O(n^{-c/2})$

▶ $\eta_k = c/k^{\beta}$ with 0 < c < 1 and $\beta \in [0, 1)$ yields the rate $O(\exp(-Cn^{(1-\beta)}))$ for some C > 0

 (Rk, ongoing work): linear rates in KL objective for constant step-size.

EMD in practice

Recall that EMD can be written at each time as:

(

$$q_{k+1}^* \propto f^{\eta_k} q_k^{*(1-\eta_k)}. \tag{2}$$

Unfortunately, running iteration (2) in practice is not feasible:

- each iteration depends on the whole density f (not some evaluations of f), maybe unknown in many applications
- ▶ (ii) even if *f* is known, computing $q_{k+1}^*(x)$ for $x \in \mathbb{R}^d$ would still be difficult due to the normalization following (2) which ensures q_{k+1}^* is a probability density function.

EMD in practice

Recall that EMD can be written at each time as:

(

$$q_{k+1}^* \propto f^{\eta_k} q_k^{*(1-\eta_k)}. \tag{2}$$

Unfortunately, running iteration (2) in practice is not feasible:

- each iteration depends on the whole density f (not some evaluations of f), maybe unknown in many applications
- (ii) even if *f* is known, computing q^{*}_{k+1}(x) for x ∈ ℝ^d would still be difficult due to the normalization following (2) which ensures q^{*}_{k+1} is a probability density function.
- \implies we propose an implementable scheme approaching (2).

Outline

Problem

Regularized Importance Sampling

Safe and Regularized Adaptive Importance Sampling

Uniform convergence of the scheme

Adaptive choice of the regularization parameter

Experiments

Safe and Regularized Adaptive Importance Sampling

We propose an *Adaptive Importance Sampling* (AIS) method which uses a sequence of proposals $(q_k)_{k\geq 0}$.

Safe and Regularized Adaptive Importance Sampling

We propose an *Adaptive Importance Sampling* (AIS) method which uses a sequence of proposals $(q_k)_{k\geq 0}$.

More specifically, as in [Delyon and Portier, 2021] we choose:

$$q_k = (1 - \lambda_k)f_k + \lambda_k q_0, \qquad \forall k \geq 1$$

- i.e. a mixture between
 - a safe density q₀ (with heavy tails compared to f), preventing too small values of q_k and high variance of IS weights,
 - a KDE estimate f_k of the target f, accelerating the convergence to f

Safe and Regularized Adaptive Importance Sampling

We propose an *Adaptive Importance Sampling* (AIS) method which uses a sequence of proposals $(q_k)_{k\geq 0}$.

More specifically, as in [Delyon and Portier, 2021] we choose:

$$q_k = (1 - \lambda_k)f_k + \lambda_k q_0, \qquad \forall k \geq 1$$

- i.e. a mixture between
 - a safe density q₀ (with heavy tails compared to f), preventing too small values of q_k and high variance of IS weights,
 - a KDE estimate f_k of the target f, accelerating the convergence to f

$$f_k(x) = \sum_{j=1}^k W_{k,j}^{(\eta_j)} \mathcal{K}_{h_k}(x-X_j), \qquad \forall x \in \mathbb{R}^d,$$

where for all $j = 1, \ldots, k$:

$$W_{k,j}^{(\eta_j)} \propto W_j^{\eta_j} = \left(\frac{f(X_j)}{q_{j-1}(X_j)}\right)^{\eta_j}, \qquad \sum_{j=1}^k W_{k,j}^{(\eta_j)} = 1.$$

Safe and Regularized Adaptive Importance sampling (SRAIS) algorithm

Inputs: The safe density q_0 , the sequences of bandwidths $(h_k)_{k=1,...,n}$, mixture weights $(\lambda_k)_{k=1,...,n}$, regularization parameters $(\eta_k)_{k=1,...,n}$.

For k = 0, 1, ..., n - 1:

1. Generate $X_{k+1} \sim q_k$.

2. Compute (a)
$$W_{k+1} = f(X_{k+1})/q_k(X_{k+1})$$

(b) $(W_{k+1,j}^{(\eta_j)})_{1 \le j \le k+1}$.

3. Return $q_{k+1} = (1 - \lambda_{k+1})f_{k+1} + \lambda_{k+1}q_0$ where $f_{k+1} = \sum_{j=1}^{k+1} W_{k+1,j}^{(\eta_j)} \mathcal{K}_{h_{k+1}}(\cdot - X_j).$

Remark: this algorithm can be used with a batch of m_k particles at each k.

SRAIS as stochastic approximation of mirror descent

Notice that

$$f_k(x) = \sum_{j=1}^k W_{k,j}^{(\eta_j)} \mathcal{K}_{h_k}(x - X_j)$$

is a stochastic approximation of the mirror descent iteration $q^*_{k+1} \propto (q^*_k)^{1-\eta_k} f^{\eta_k}$. Indeed,

$$\mathbb{E}_{X_{j}\sim q_{j-1}}[W_{j}^{\eta_{j}}K_{h_{k}}(x-X_{j})]=(f^{\eta_{j}}q_{j-1}^{1-\eta_{j}}\star K_{h_{k}})(x),$$

which approximates $f^{\eta_j}q_{j-1}^{1-\eta_j}$ when the bandwidth h_k is small.

Outline

Problem

Regularized Importance Sampling

Safe and Regularized Adaptive Importance Sampling

Uniform convergence of the scheme

Adaptive choice of the regularization parameter

Experiments

Convergence of the scheme?

Does the kernel density estimate (KDE) f_k converge to the target distribution f?

Recall the hyperparameters of the algorithm :

- the safe density q₀ in the proposal (mixture between q₀ and f_k)
- mixture weights $(\lambda_k)_{k\geq 0}$
- ▶ KDE bandwidths $(h_k)_{k\geq 0}$
- regularization parameters $(\eta_k)_{k\geq 0}$

(A1)

1. The sequence $(\lambda_k)_{k\geq 1}$ is valued in (0, 1], nonincreasing, and $\lim_{k\to\infty} \lambda_k = 0$ and $\lim_{k\to\infty} \log(k)/(k\lambda_k) = 0$.

(A1)

- 1. The sequence $(\lambda_k)_{k\geq 1}$ is valued in (0, 1], nonincreasing, and $\lim_{k\to\infty} \lambda_k = 0$ and $\lim_{k\to\infty} \log(k)/(k\lambda_k) = 0$.
- 2. The sequence $(h_k)_{k\geq 1}$ is valued in \mathbb{R}^+ , nonincreasing, and $\lim_{k\to\infty} h_k = 0$ and $\lim_{k\to\infty} \log(k)/(kh_k^d \lambda_k) = 0$.

(A1)

- 1. The sequence $(\lambda_k)_{k\geq 1}$ is valued in (0, 1], nonincreasing, and $\lim_{k\to\infty} \lambda_k = 0$ and $\lim_{k\to\infty} \log(k)/(k\lambda_k) = 0$.
- 2. The sequence $(h_k)_{k\geq 1}$ is valued in \mathbb{R}^+ , nonincreasing, and $\lim_{k\to\infty} h_k = 0$ and $\lim_{k\to\infty} \log(k)/(kh_k^d \lambda_k) = 0$.
- 3. The sequence $(\eta_k)_{k\geq 1}$ is valued in (0, 1], and $\lim_{k\to\infty} \eta_k = 1$, $\lim_{k\to\infty} (1-\eta_k) \log(h_k) = 0$ and $\lim_{k\to\infty} (1-\eta_k) \log(\lambda_{k-1}) = 0$.

(A1)

- 1. The sequence $(\lambda_k)_{k\geq 1}$ is valued in (0, 1], nonincreasing, and $\lim_{k\to\infty} \lambda_k = 0$ and $\lim_{k\to\infty} \log(k)/(k\lambda_k) = 0$.
- 2. The sequence $(h_k)_{k\geq 1}$ is valued in \mathbb{R}^+ , nonincreasing, and $\lim_{k\to\infty} h_k = 0$ and $\lim_{k\to\infty} \log(k)/(kh_k^d \lambda_k) = 0$.
- 3. The sequence $(\eta_k)_{k\geq 1}$ is valued in (0, 1], and $\lim_{k\to\infty} \eta_k = 1$, $\lim_{k\to\infty} (1-\eta_k) \log(h_k) = 0$ and $\lim_{k\to\infty} (1-\eta_k) \log(\lambda_{k-1}) = 0$.

(A2) The density q_0 is bounded and there exists c > 0 such that for all $x \in \mathbb{R}^d$, $q_0(x) \ge cf(x)$.

(A1)

- 1. The sequence $(\lambda_k)_{k\geq 1}$ is valued in (0, 1], nonincreasing, and $\lim_{k\to\infty} \lambda_k = 0$ and $\lim_{k\to\infty} \log(k)/(k\lambda_k) = 0$.
- 2. The sequence $(h_k)_{k\geq 1}$ is valued in \mathbb{R}^+ , nonincreasing, and $\lim_{k\to\infty} h_k = 0$ and $\lim_{k\to\infty} \log(k)/(kh_k^d \lambda_k) = 0$.
- 3. The sequence $(\eta_k)_{k\geq 1}$ is valued in (0, 1], and $\lim_{k\to\infty} \eta_k = 1$, $\lim_{k\to\infty} (1-\eta_k) \log(h_k) = 0$ and $\lim_{k\to\infty} (1-\eta_k) \log(\lambda_{k-1}) = 0$.

(A2) The density q_0 is bounded and there exists c > 0 such that for all $x \in \mathbb{R}^d$, $q_0(x) \ge cf(x)$.

(A3) The function *f* is nonnegative, bounded by $B_f \ge 0$ and is I_f -Lipschitz.

(A1)

- 1. The sequence $(\lambda_k)_{k\geq 1}$ is valued in (0, 1], nonincreasing, and $\lim_{k\to\infty} \lambda_k = 0$ and $\lim_{k\to\infty} \log(k)/(k\lambda_k) = 0$.
- 2. The sequence $(h_k)_{k\geq 1}$ is valued in \mathbb{R}^+ , nonincreasing, and $\lim_{k\to\infty} h_k = 0$ and $\lim_{k\to\infty} \log(k)/(kh_k^d \lambda_k) = 0$.
- 3. The sequence $(\eta_k)_{k\geq 1}$ is valued in (0, 1], and $\lim_{k\to\infty} \eta_k = 1$, $\lim_{k\to\infty} (1-\eta_k) \log(h_k) = 0$ and $\lim_{k\to\infty} (1-\eta_k) \log(\lambda_{k-1}) = 0$.

(A2) The density q_0 is bounded and there exists c > 0 such that for all $x \in \mathbb{R}^d$, $q_0(x) \ge cf(x)$.

(A3) The function *f* is nonnegative, bounded by $B_f \ge 0$ and is I_f -Lipschitz.

(A4) The kernel *K* is bounded by $K_{\infty} \ge 0$ and is L_K -Lipschitz with $L_K > 0$. Moreover, $\int K(u) du = 1$, $\int ||u|| K(u) du < \infty$, $\int K^{1/2}(u) d(u) < \infty$ and $\int ||u|| K(u)^{1/2} du < \infty$.

Uniform convergence of the scheme

Proposition: Assume **A1-A4**. Then, for any r > 0:

$$\sup_{\|x\|\leq k^r} |f_k(x)-f(x)| o 0$$
 as $k o\infty$ a.s.

Proof: $f_k = N_k/D_k$ where

$$N_k(x) = rac{1}{k} \sum_{j=1}^k W_j^{\eta_j} K_{h_k}(x - X_j), \quad D_k = rac{1}{k} \sum_{j=1}^k W_j^{\eta_j}.$$

Uniform convergence of the scheme

Proposition: Assume A1-A4. Then, for any r > 0:

$$\sup_{\|x\|\leq k^r} |f_k(x)-f(x)| o 0$$
 as $k o\infty$ a.s.

Proof: $f_k = N_k/D_k$ where

$$N_{k}(x) = \frac{1}{k} \sum_{j=1}^{k} W_{j}^{\eta_{j}} K_{h_{k}}(x - X_{j}), \quad D_{k} = \frac{1}{k} \sum_{j=1}^{k} W_{j}^{\eta_{j}}.$$

$$N_{k} - f = \underbrace{\frac{1}{k} \sum_{j=1}^{k} \left\{ W_{j}^{\eta_{j}} K_{h_{k}}(x - X_{j}) - \{f^{\eta_{j}} q_{j-1}^{1-\eta_{j}}\} \star K_{h_{k}}(x)\}}_{(1)} + \underbrace{\left\{ \left(\frac{1}{k} \sum_{j=1}^{k} (f^{\eta_{j}} q_{j-1}^{1-\eta_{j}} - f)\right) \star K_{h_{n}} \right\}}_{(2)} + \underbrace{\{f \star K_{h_{k}} - f\}}_{(3)}.$$

(1): avg of martingale increments, (2): regularization bias (\rightarrow 0 as $\eta_k \rightarrow$ 1), (3) KDE bias (\rightarrow 0 as $h_k \rightarrow$ 0).

21/32

Outline

Problem

Regularized Importance Sampling

Safe and Regularized Adaptive Importance Sampling

Uniform convergence of the scheme

Adaptive choice of the regularization parameter

Experiments

Adaptive Choice of Regularization (RAR)

Our conditions for uniform convergence require that the sequence $(\eta_k)_{k\geq 1}$ converges to 1. We propose an adaptive way to construct it.

Adaptive Choice of Regularization (RAR)

Our conditions for uniform convergence require that the sequence $(\eta_k)_{k\geq 1}$ converges to 1. We propose an adaptive way to construct it.

Idea: Draw m_k i.i.d samples $X_{k,1}, \ldots, X_{k,m_k}$ from q_{k-1} . Let $\mathbb{P} = \sum_{l=1}^{m_k} W_{k,l} \delta_{X_{k,l}}$ and $\mathbb{Q} = \sum_{l=1}^{m_k} \frac{1}{m_k} \delta_{X_{k,l}}$ the reweighted and uniform distribution on the particles. \implies If $q_{k-1} = f$, IS weights = 1 and $\mathbb{P} = \mathbb{Q}$. \implies **penalize the divergence between** \mathbb{P} and \mathbb{Q} !

Adaptive Choice of Regularization (RAR)

Our conditions for uniform convergence require that the sequence $(\eta_k)_{k\geq 1}$ converges to 1. We propose an adaptive way to construct it.

Idea: Draw m_k i.i.d samples $X_{k,1}, \ldots, X_{k,m_k}$ from q_{k-1} . Let $\mathbb{P} = \sum_{l=1}^{m_k} W_{k,l} \delta_{X_{k,l}}$ and $\mathbb{Q} = \sum_{l=1}^{m_k} \frac{1}{m_k} \delta_{X_{k,l}}$ the reweighted and uniform distribution on the particles. \implies If $q_{k-1} = f$, IS weights = 1 and $\mathbb{P} = \mathbb{Q}$. \implies penalize the divergence between \mathbb{P} and \mathbb{Q} !

We propose to use Renyi's α -divergences and set:

$$\eta_{k,\alpha} = 1 - \frac{D_{\alpha}(\mathbb{P}||\mathbb{Q})}{\log(m_k)}, \ D_{\alpha}(\mathbb{P}||\mathbb{Q}) = \frac{1}{\alpha - 1} \log \left(\sum_{\ell=1}^{m_k} W_{k,\ell}^{\alpha} m_k^{\alpha - 1} \right)$$

$$\eta_{k,\alpha} = 1 - \frac{D_{\alpha}(\mathbb{P}||\mathbb{Q})}{\log(m_k)}, \ D_{\alpha}(\mathbb{P}||\mathbb{Q}) = \frac{1}{\alpha - 1} \log\left(\sum_{\ell=1}^{m_k} W_{k,\ell}^{\alpha} m_k^{\alpha - 1}\right).$$
(3)

Proposition: Let $\alpha \in [0, 1]$ and let $(\eta_{k,\alpha})_{k \ge 1}$ be the sequence defined by (3) for all $k \ge 1$. Then, we have:

- 1. The sequence $(\eta_{k,\alpha})_{k\geq 1}$ is valued in [0, 1], with $\eta_{k,\alpha} = 1$ iff $\mathbb{P} = \mathbb{Q}$;
- **2**. $0 \le \eta_{k,1} \le \eta_{k,\alpha} \le 1$;
- 3. Further assume that $(q_k)_{k\geq 1}$ is a sequence of probability density functions s.t. $\lim_{k\to\infty} |q_k(x) f(x)| = 0$ almost everywhere and that $m_k = m$ for all $k \geq 1$ (fixed batch size). Then, $\lim_{k\to\infty} \eta_{k,\alpha} = 1$ in L_1 .

Renyi's $\alpha\text{-divergences}$

α	Definition	Notes
$\alpha \rightarrow 1$	$\int p(\boldsymbol{\theta}) \log \frac{p(\boldsymbol{\theta})}{d\boldsymbol{\theta}} d\boldsymbol{\theta}$	Kullback-Leibler (KL) divergence,
u / 1	$\int P(0) \log q(\mathbf{\theta}) d0$	used in VI (KL[$q p]$) and EP (KL[$p q]$)
lpha=0.5	$-2\log(1-\mathrm{Hel}^2[p q])$	function of the square Hellinger distance
$\alpha \rightarrow 0$	$-\log\int_{p(oldsymbol{ heta})>0}q(oldsymbol{ heta})doldsymbol{ heta}$	zero when $\operatorname{supp}(q) \subseteq \operatorname{supp}(p)$
		(not a divergence)
$\alpha = 2$	$-\log(1-\chi^2[p q])$	proportional to the χ^2 -divergence
$\alpha \to +\infty$	$\log \max_{\boldsymbol{\theta} \in \Theta} \frac{p(\boldsymbol{\theta})}{q(\boldsymbol{\theta})}$	worst-case regret in
		minimum description length principle [24]

Table 1: Special cases in the Rényi divergence family.

Outline

Problem

Regularized Importance Sampling

Safe and Regularized Adaptive Importance Sampling

Uniform convergence of the scheme

Adaptive choice of the regularization parameter

Experiments

Toy Experiments



Figure: Logarithm of the average squared error for SRAIS for constant values of η or Adaptive η , over 50 replicates. 4×10^4 particles sampled from initial density, then $m_k = 18 \times 10^3$ particles from q_k at each $k \ge 1$.

Different target densities ($\phi_{\Sigma} = \mathcal{N}(\mathbf{0}_{d}, \Sigma)$), initial densities have different means/variance than the target:

- "Cold Start" $f_1(x) = \phi_{\Sigma}(x 5\mathbf{1}_d/\sqrt{d}), \Sigma = (0.16/d)\mathbf{I}_d$
- "Gaussian Mixture" $f_2(x) = 0.5\phi_{\Sigma}(x - \mathbf{1}_d/(2\sqrt{d})) + 0.5\phi_{\Sigma}(x + \mathbf{1}_d/(2\sqrt{d}))$

• "Anisotropic Gaussian Mixture" $f_3(x) = 0.25\phi_V(x - \mathbf{1}_d/(2\sqrt{d})) + 0.75\phi_V(x + \mathbf{1}_d/(2\sqrt{d})),$ $V = (.4/\sqrt{d})^2 \operatorname{diag}(10, 1, ..., 1)$

Evolution of Adaptive Regularization



Figure: Boxplot of the values of $(\eta_{k,\alpha})_{k\geq 1}$ obtained from RAR (Adaptive η), with $\alpha = 0.5$.

- at the beginning of the algorithm when the policy is poor, the value of η_k is automatically set to a small value (leading to a uniformization of the weights)
- when the policy becomes better the value of η_{k,α} converges to 1.

Bayesian Logistic Regression (Waveform dataset, 5000 datapoints in d = 22)



Figure: Left plot: Average accuracy over 100 trials of different learning policies $(\eta_{k,\alpha})_{k\geq 1}$ for Bayesian Logistic Regression on the Waveform dataset. Right plot: Averaged values of the learning policy $(\eta_{k,\alpha})_{k\geq 1}$ associated to each choice of α .

- a proper tuning of the parameter α allows us to outperform (η_k)_{k≥1} constant and equal to 1
- the case α = 0.2 yielding the best results here overall in terms of speed and accuracy

Contributions:

- We proposed a new algorithm for Adaptive Importance Sampling, that regularizes the importance weights by raising them to a certain power
- This algorithm is related to mirror descent on the space of probability distributions
- It enjoys a uniform convergence guarantee under mild assumptions on the target, safe density, and hyperparameters
- It outperforms numerically constant values of η

Future work:

Non-asymptotic analysis of the scheme

Future work:

- Non-asymptotic analysis of the scheme
- Adaptive schedules for other hyperparameters

Future work:

- Non-asymptotic analysis of the scheme
- Adaptive schedules for other hyperparameters
- Replace KDE by a more scalable approximate ? e.g. normalizing flows [Papamakarios et al., 2021]:
 - generative models that transform a simple $q \in \mathcal{P}(\mathbb{R}^d)$ into $\tilde{f} \in \mathcal{P}(\mathbb{R}^d)$ ($\tilde{f} \approx f$)

through a sequence of invertible transformations:

$$ilde{f} = G_{\#}q = (g_1 \circ \ldots g_n)_{\#}q, \quad g_i: \mathbb{R}^d o \mathbb{R}^d$$

 $\blacktriangleright \implies$ sampling and evaluating the density are efficient

$$ilde{f}(y) = q(z_0) \prod_{i=1}^n \left| det\left(rac{\partial g_i}{\partial z_{i-1}}
ight) \right|^{-1}, \quad y, z_0 \in \mathbb{R}^d$$

where $\frac{\partial g}{\partial z_{i-1}}$ is the Jacobian matrix of g_i .

Future work:

- Non-asymptotic analysis of the scheme
- Adaptive schedules for other hyperparameters
- Replace KDE by a more scalable approximate ? e.g. normalizing flows [Papamakarios et al., 2021]:
 - generative models that transform a simple $q \in \mathcal{P}(\mathbb{R}^d)$ into $\tilde{f} \in \mathcal{P}(\mathbb{R}^d)$ ($\tilde{f} \approx f$)

through a sequence of invertible transformations:

$$ilde{f} = G_{\#}q = (g_1 \circ \ldots g_n)_{\#}q, \quad g_i: \mathbb{R}^d o \mathbb{R}^d$$

 $\blacktriangleright \implies$ sampling and evaluating the density are efficient

$$ilde{f}(y) = q(z_0) \prod_{i=1}^n \left| det\left(rac{\partial g_i}{\partial z_{i-1}}
ight) \right|^{-1}, \quad y, z_0 \in \mathbb{R}^d$$

where $\frac{\partial g}{\partial z_{i-1}}$ is the Jacobian matrix of g_i . Thank you !

References I

- Delyon, B. and Portier, F. (2021).
 Safe adaptive importance sampling: A mixture approach. *The Annals of Statistics*, 49(2):885–917.
- Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. (2021).
 Normalizing flows for probabilistic modeling and inference. *Journal of Machine Learning Research*, 22(57):1–64.