

ITO DIFFUSION APPROXIMATION OF UNIVERSAL ITO CHAINS FOR SAMPLING, OPTIMIZATION AND BOOSTING

Aleksei Ustimenko

ShareChat, UK

aleksei.ustimenko@sharechat.co

Aleksandr Beznosikov

Innopolis University, Skoltech, RANEPa, Yandex

anbeznosikov@gmail.com

Aleksei Ustimenko, International Conference on Learning Representations 2024

What do we study?

- $X_{k+1} = X_k + \eta(b(X_k) + \delta_k) + \sqrt{\eta^{1+\gamma}}(\sigma(X_k) + \Delta_k)\epsilon_k(X_k),$
- On the first glance — a generic Markov chain with a drift function $b(x)$ and a covariance $\sigma(x)$
 - Assume those functions are Lipschitz-smooth, not even dissipative
- δ_k, Δ_k — a bias/covariance shift of an order $O(\eta^\alpha)$
- $\gamma \in \{0,1\}$ to cover SGD-like chains in the same analysis
- We want to bound Wasserstein-2 distance to the Itô-process:
 - $dZ_t = b(Z_t)dt + \sqrt{\eta^\gamma}\sigma(Z_t)dW_t.$

What do we study?

- $X_{k+1} = X_k + \eta(b(X_k) + \delta_k) + \sqrt{\eta^{1+\gamma}}(\sigma(X_k) + \Delta_k)\epsilon_k(X_k),$
- The noise-function $\epsilon_k(x)$ — different for every step k and covers non-gaussian state-dependent isotropic noise, possibly non-gaussian
- Assume it satisfies a weak version of the CLT in Wasserstein-2 metric
$$\mathbb{E}_\epsilon \epsilon_k(x) = 0_d, \quad \mathbb{E}_\epsilon \epsilon_k \epsilon_k^\top(x) = I_d, \quad \mathcal{W}_2^2\left(\mathcal{N}(0_d, SI_d), \mathcal{L}\left(\sum_{k=0}^{S-1} \epsilon_k(x)\right)\right) \leq M_\epsilon^2 \eta^\beta.$$
- Example: any noise of the form $\Sigma(x)\xi, \mathbb{E}\|\xi\|^4 < \infty$

Why such a generic form is of interest?

Table 1: Matching different methods and frameworks with the parameters of the Ito chain equation 1 and Assumption 1.

| | Case | γ | α | β | $b(X_k)$ | δ_k | $\sigma(X_k)$ | Δ_k |
|--------------|---|----------|----------------|----------|---|--|---------------------------------------|---|
| Dynamics | GLD | 0 | $\infty^{(1)}$ | ∞ | $-\nabla_x f(X_k)$ | 0 | $\sqrt{\frac{2}{\tau}} I_d^{(2)}$ | 0 |
| | SGLD (Gelfand et al., 1992) | 0 | $\frac{1}{2}$ | 1 | $-\nabla_x f(X_k)$ | 0 | $\sqrt{\frac{2}{\tau}} I_d$ | $\sqrt{\frac{2}{\tau}} I_d + \eta \text{Cov}(\widehat{\nabla}) - \sqrt{\frac{2}{\tau}} I_d$ |
| | SGLD with smoothing (Chatterji et al., 2020) | 0 | $\frac{1}{2}$ | 1 | $-\nabla_x f(X_k)$ | $\nabla_x (f(X_k) - \mathbb{E}_\varepsilon f(X_k + \eta^{\frac{1}{2}} \varepsilon))$ | $\sqrt{\frac{2}{\tau}} I_d$ | $\sqrt{\frac{2}{\tau}} I_d + \eta \text{Cov}(\widehat{\nabla}) - \sqrt{\frac{2}{\tau}} I_d$ |
| Optimization | SGD (Robbins & Monro, 1951) | 1 | ∞ | 0 | $-\nabla_x f(X_k)$ | 0 | $\sqrt{\text{Cov}(\nabla)}$ | 0 |
| | SGDA (Dem'yanov & Pevnyi, 1972) | 1 | ∞ | 0 | $(-\nabla_x f(X_k, Y_k), \nabla_y f(X_k, Y_k))$ | 0 | $\sqrt{\text{Cov}(\nabla)}$ | 0 |
| | SA-FP (Bailion et al., 1978) | 1 | ∞ | 0 | $F(X_k) - X_k$ | 0 | $\sqrt{\text{Cov}(\widehat{F})}$ | 0 |
| | SA (Bailion et al., 1978) | 1 | ∞ | 0 | $H(X_k) - a^{(3)}$ | 0 | $\sqrt{\text{Cov}(\widehat{H})}$ | 0 |
| Boosting | SGB (Friedman, 2001) | 1 | ∞ | 0 | $-P(X_k) \nabla_x f(X_k)$ | 0 | $\sqrt{\text{Cov}(\widehat{\nabla})}$ | 0 |
| | SGLB (Ustimenko & Prokhorenkova, 2021) ⁽⁴⁾ | 0 | $\frac{1}{2}$ | 0 | $-P(X_k) \nabla_x f(X_k)$ | 0 | $\sqrt{\frac{2}{\tau}} I_d$ | $\sqrt{\eta \text{Cov}(\widehat{\nabla})}$ |
| | SGLB-O (Ustimenko & Prokhorenkova, 2021) ⁽⁵⁾ | 0 | $\frac{1}{4}$ | 0 | $-P_\infty \nabla_x f(X_k)$ | $(P_\infty - P(X_k)) \nabla_x f(X_k)$ | $\sqrt{\frac{2}{\tau}} I_d$ | $\sqrt{\eta \text{Cov}(\widehat{\nabla})}$ |

⁽¹⁾ η^∞ means that the terms multiplied by it vanish, i.e., we can take α as large as we desire when calculating overall approximation error.

⁽²⁾ τ refers to inverse diffusion temperature.

⁽³⁾ a is any constant. Stochastic Approximation tries to solve $H(x) = a$.

⁽⁴⁾ SGLB here is defined as in the original paper, but here we ignore smoothing applied to the trees selection algorithm.

⁽⁵⁾ "O" stands for "original", i.e., as presented in the original paper. In that case, such coefficients appear if we take the distribution of trees as in the paper.

How it compares with prior works?

Table 2: Comparison of the theoretical setups and results on Markov chains and diffusions analysis.

| Reference | Noise | | Generator, i.e. $b(\cdot)$ | | | |
|-----------------------------------|--|------------|----------------------------|-----------------|------------------------|-----------------|
| | Distribution | Dependence | Non-convex | Non-dissipative | Non-uniformly elliptic | \mathcal{W}_2 |
| (Raginsky et al., 2017) | $\mathcal{N}+\text{SG}$ | ✓ | ✓ | ✗ | ✓ | ✓ |
| (Dalalyan, 2017) | \mathcal{N} | ✗ | ✗ | ✗ | ✗ | ✓ |
| (Cheng et al., 2018) | \mathcal{N} | ✗ | ✓ | ✗ | ✗ | ✗ |
| (Erdogdu et al., 2018) | $\mathcal{N}+\text{SG}$ | ✓ | ✓ | ✗ | ✗ | ✗ |
| (Durmus & Moulines, 2019) | \mathcal{N} | ✗ | ✗ | ✗ | ✗ | ✓ |
| (Ma et al., 2019) | \mathcal{N} | ✗ | ✓ | ✗ | ✗ | ✗ |
| (Li et al., 2019b) | \mathcal{N} | ✓ | ✗ | ✗ | ✓ | ✓ |
| (Chatterji et al., 2020) | \mathcal{N} | ✗ | ✗ | ✗ | ✗ | ✓ |
| (Feng et al., 2019) | $\ \epsilon\ \leq \text{const a.s.}$ | ✗ | ✓ | ✗ | ✗ | ✗ |
| (Orvieto & Lucchi, 2018) | \mathcal{N} | ✗ | ✓ | ✗ | ✗ | ✗ |
| (Ankirchner & Perko, 2021) | \mathcal{N} | ✗ | ✓ | ✗ | ✓ | ✗ |
| (Hu et al., 2017) | \mathcal{N} | ✗ | ✓ | ✗ | ✗ | ✗ |
| (Xie et al., 2021) | \mathcal{N} | ✗ | ✗ | ✗ | ✗ | ✗ |
| (Ustimenko & Prokhorenkova, 2021) | Mixture \mathcal{N} | ✓ | ✓ | ✗ | ✓ | ✗ |
| (Cheng et al., 2020) | $\mathcal{N}+\text{SG}$ | ✓ | ✓ | ✗ | ✗ | ✗ |
| (Li et al., 2019a) | $\mathbb{E}\ \epsilon\ ^4 \leq \text{const}$ | ✓ | ✓ | ✗ | ✓ | ✗ |
| Ours | $\mathbb{E}\ \epsilon\ ^4 \leq \text{const}$ | ✓ | ✓ | ✓ | ✗ | ✓ |

Table 3: Comparison of guarantees on discretization error in our work with the literature.

| | Noise | Reference | Rate |
|------|----------|-------------------------------------|-----------------------------------|
| SGLD | Gaussian | (Muzellec et al., 2020) | $\mathcal{O}(\eta^{\frac{1}{4}})$ |
| | | Ours | $\mathcal{O}(\eta^{\frac{1}{4}})$ |
| | any | Ours | $\mathcal{O}(\eta^{\frac{1}{4}})$ |
| SGD | Gaussian | (Cheng et al., 2020) ⁽¹⁾ | $\mathcal{O}(\eta^{\frac{1}{8}})$ |
| | | Ours | $\mathcal{O}(\eta^{\frac{3}{4}})$ |
| | any | (Cheng et al., 2020) ⁽¹⁾ | $\mathcal{O}(\eta^{\frac{1}{8}})$ |
| | | Ours | $\mathcal{O}(\eta^{\frac{1}{2}})$ |
| SGB | any | Ours | $\mathcal{O}(\eta^{\frac{1}{2}})$ |
| SGLB | any | Ours | $\mathcal{O}(\eta^{\frac{1}{4}})$ |

⁽¹⁾ \mathcal{W}_1 - distance

High level idea

- optimal transport $(\sqrt{S}\zeta_k^S(x), \epsilon_k^S(x)) \sim \Pi_*(\mathcal{N}(0_d, SI_d), \mathcal{L}(\epsilon_k^S(x)))$,
- chain with Gaussian noise $Y_{\bar{\eta}(k+1)}^X = Y_{\bar{\eta}k}^X + \bar{\eta}b(Y_{\bar{\eta}k}^X) + \underbrace{\sqrt{\bar{\eta}\eta^\gamma} \sigma(Y_{\bar{\eta}k}^X) \zeta_k^S(X_{S_k})}_{\text{coupling via noise}}$,
 - **it won't work as we don't assume dissipativity of $b(x)$!**
- “regularize it”: $Y_{\bar{\eta}(k+1)}^X = Y_{\bar{\eta}k}^X + \bar{\eta}b(Y_{\bar{\eta}k}^X) + \sqrt{\bar{\eta}\eta^\gamma} \sigma(Y_{\bar{\eta}k}^X) \zeta_k^S(X_{S_k}) - \underbrace{L\bar{\eta}(Y_{\bar{\eta}k}^X - X_{S_k})}_{\text{Window coupling}}$.
- Interpolate $dY_t = (b(Y_t) + (g_t^* + g_t^S))dt + \sqrt{\eta^\gamma}(\sigma(Y_t) + \Sigma_t^*)dW_t$,
- Yet another process to prepare for the Girsanov theorem:

$$dZ_t^Y = (b(Z_t^Y) + g_t^S)dt + \underbrace{L_1(Y_t - Z_t^Y)}_{\text{Window coupling}}dt + \sqrt{\eta^\gamma} \sigma(Z_t^Y) dW_t,$$

- Use a weaker version of the Girsanov theorem, which we proved in the work

Main results

Theorem 1 (One-time Girsanov formula for mixed Ito/adapted coefficients). *Assume that $(G_t)_{t \geq 0}$ is a $(W_t)_{t \geq 0}$ -adapted process with an integrable by $t \geq 0$ second moment. Consider two SDEs run using two independent Brownian processes:*

$$\begin{aligned}dZ_t &= b(Z_t)dt + G_t dt + \sigma(Z_t)dW_t, \\dZ_t^* &= b(Z_t^*)dt + \sigma(Z_t^*)d\widetilde{W}_t.\end{aligned}$$

Let $\sigma_0 > 0$ be the minimal possible eigenvalue of $\sigma(x)$. Then we have that for any time horizon $T \geq 0$ the following bound holds:

$$D_{\text{KL}}(\mathcal{L}(Z_T) \parallel \mathcal{L}(Z_T^*)) \leq \sigma_0^{-2} \int_0^T \mathbb{E} \|G_t\|^2 dt < \infty.$$

Theorem 2. *Let Assumption [1](#) hold. Then for all $k \in \mathbb{N}_0$:*

$$\mathcal{W}_2(\mathcal{L}(X_{k'}), \mathcal{L}(Z_{k'\eta})) = \mathcal{O}\left(\left(1 + (k'\eta)^{\frac{1}{2}}\right)e^{\mathcal{O}(k'\eta)}\eta^\theta + (k'\eta)^{\frac{1}{4}}e^{\mathcal{O}(k'\eta)}\eta^{\frac{\theta}{2} + \frac{\gamma}{4}}\right),$$

where $\theta = \min\left\{\alpha; \frac{(\gamma+1)(1+\chi_0) + (\gamma+\beta)(1-\chi_0)}{4}\right\}$.

Corollary 5 (SGD with Gaussian noise). *Under the conditions of Theorem [2](#), if $\chi_0 = 1$, $\gamma = 1$, $\alpha \geq 1$, then $\theta = 1$ and $\mathcal{W}_2(\mathcal{L}(X_k), \mathcal{L}(Z_{k\eta})) = \mathcal{O}\left(\eta^{\frac{3}{4}}\right)$.*

Thank you