

# **Asymptotic Behaviors of Projected Stochastic Approximation: A Jump Diffusion Perspective**

**Jiadong Liang; Yuze Han; Xiang Li; Zhihua Zhang**

# Loopless Projection Stochastic Approximation

- We aim to solve the following problem

$$\min_{\mathbf{x}} \mathbb{E}_{\zeta \sim \mathcal{D}} f(\mathbf{x}, \zeta) \text{ subject to } \mathbf{A}^\top \mathbf{x} = \mathbf{0}$$

- The LPSA first performs as  $\mathbf{x}_{n+\frac{1}{2}} = \mathbf{x}_n - \eta_n \nabla f(\mathbf{x}_n) + \eta_n \xi_n$  with a martingale difference sequence  $\{\xi_n\}$ .
- Then we independently cast a coin with the head probability  $p_n$  and obtain  $\omega_n \sim \text{Bernoulli}(p_n)$ . If  $\omega_n = 1$ , we perform one step of projection onto the null space of  $\mathbf{A}^\top$ . Otherwise, we let  $\mathbf{x}_{n+1} = \mathbf{x}_{n+\frac{1}{2}}$ .
- It's obvious that Local SGD is a specialized case of LPSA under federated learning scenario.

# Convergence Rate Analysis

- Let  $\eta_n = \eta_0 n^{-\alpha}$ ;  $p_n = \min\{\eta_n^\beta, 1\}$  with  $\beta \in [0, 1)$ .

- **Theorem 3.1**

Under appropriate assumptions, for (i)  $0 < \alpha < 1$  or (ii)  $\alpha = 1$  with  $\eta_0 > 2/\mu$  ( $\mu$  is the strong convexity parameter), we have

$$\mathbb{E} \left\| \mathbf{u}_n - \mathbf{x}^\star \right\|^2 = \mathcal{O} \left( n^{-\alpha \min\{1, 2-2\beta\}} \right)$$

Where  $\mathbf{u}_n = \mathcal{P}_{A^\perp}(\mathbf{x}_n)$ .

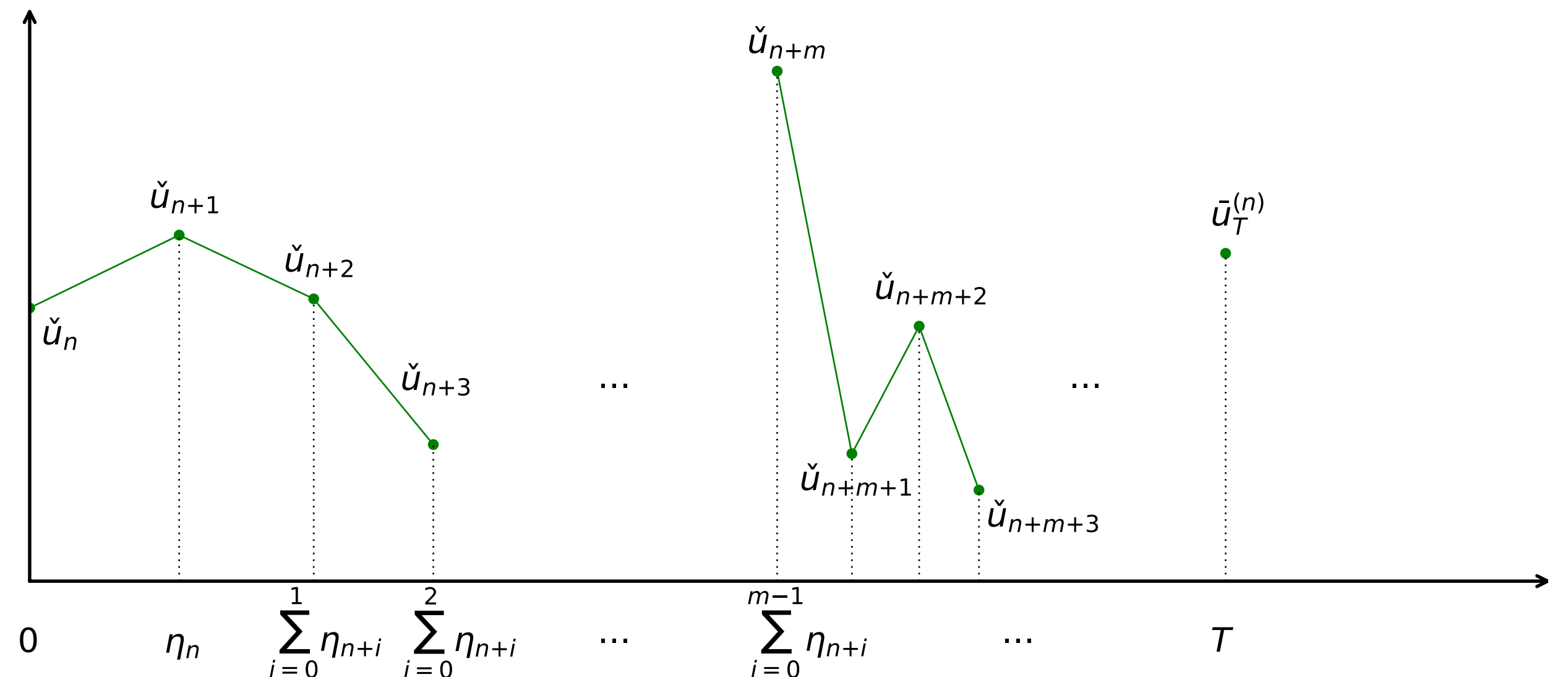
- As  $\beta$  decreases, i.e. the projection happens more frequently,  $\mathbb{E} \left\| \mathbf{u}_n - \mathbf{x}^\star \right\|^2$  converges faster. What's more, we can find a phase transition when  $\beta$  goes cross 0.5.

# Asymptotic Behavior via Diffusion Approximation

- Frequent Projection  $\beta \in [0, 1/2)$

- Let  $\check{u}_n := \frac{u_n - x^*}{\sqrt{\eta_{n-1}}}$ . And let  $\bar{u}_t^{(n)}$  be the continuous random process which starts at  $\check{u}_n$  and takes value  $\check{u}_{n+m}$  at time point  $\eta_n + \dots + \eta_{n+m-1}$ .

- The trajectory is presented in the form shown on the right.



# Asymptotic Behavior via Diffusion Approximation

- **Theorem 3.3.**

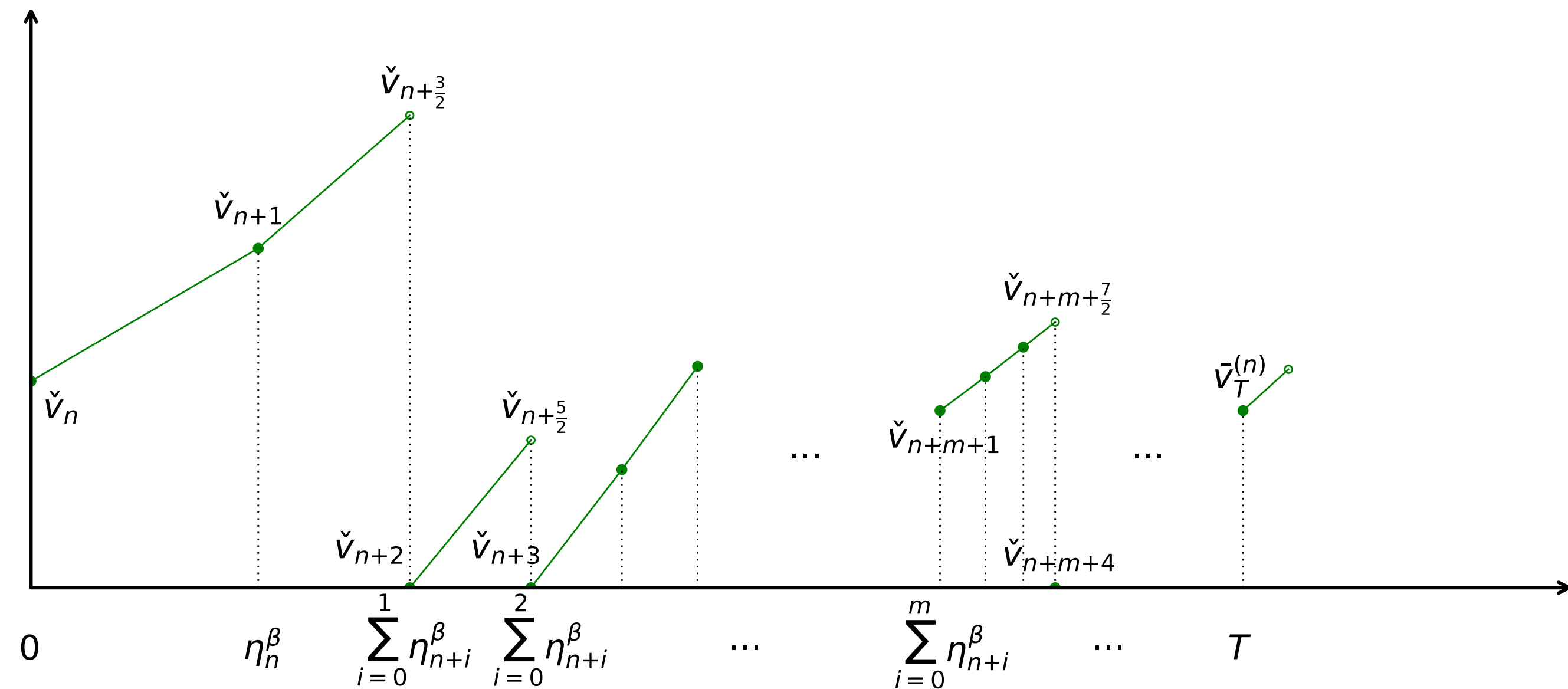
Let regular assumptions hold. Then the sequence of random processes  $\{\bar{\mathbf{u}}_t^{(n)} : t \geq 0\}_{n=1}^{\infty}$  converges weakly to the stationary weak solution of the following SDE:

$$d\mathbf{X}_t = -\mathcal{P}_{A^\perp} \left( \nabla^2 f(\mathbf{x}^\star) - \frac{1}{2\eta_0} \mathbf{1}_{\{\alpha=1\}} \mathbf{I}_d \right) \mathbf{X}_t dt + \mathcal{P}_{A^\perp} \Sigma(\mathbf{x}^\star)^{\frac{1}{2}} d\mathbf{W}_t.$$

Further, the rescaled sequence  $\{\check{\mathbf{u}}_n\}_{n=1}^{\infty}$  converges weakly to the invariant distribution of this dynamics.

# Asymptotic Behavior via Jump Approximation

- Occasional Projection  $\beta \in (1/2, 1)$ .
- Let  $\mathbf{v}_n = \mathcal{P}_A(\mathbf{x}_n)$  and  $\check{\mathbf{v}}_n = \eta_{n-1}^{\beta-1} \mathbf{v}_n$ . And we define  $\bar{\mathbf{v}}_t^{(n)}$  as the cadlag process which starts at  $\check{\mathbf{v}}_n$  and take values  $\check{\mathbf{v}}_{n+m+\frac{1}{2}}, \check{\mathbf{v}}_{n+m+1}$  at time points  $(\eta_n^\beta + \dots + \eta_{n+m-1}^\beta)^-$  and  $\eta_n^\beta + \dots + \eta_{n+m-1}^\beta$  respectively.



# Asymptotic Behavior via Jump Approximation

- **Theorem 3.4**

Let regular assumptions hold. Then the sequence of cadlag stochastic processes  $\{\bar{\mathbf{v}}_t^{(n)} : t \geq 0\}_{n=1}^{\infty}$  weakly converges to the stationary weak solution of the following Jump-SDE:

$$d\mathbf{Y}_t = -\nabla f(\mathbf{x}^*) dt - \mathbf{Y}_{t-} \cdot \mathbf{N}_\gamma(dt).$$

Further, the rescaled sequence  $\{\check{\mathbf{v}}_n\}_{n=1}^{\infty}$  weakly converges to the invariant

distribution of this dynamics, i.e.,  $\frac{\nabla f(\mathbf{x}^*)}{\|\nabla f(\mathbf{x}^*)\|} \cdot \mathcal{E}\left(\frac{\|\nabla f(\mathbf{x}^*)\|}{\gamma}\right)$

# Asymptotic Behavior via Jump Approximation

- Corollary 1

Let regular assumptions hold. Then for  $\beta \in (1/2, 1)$ ,

$\hat{\mathbf{u}}_n := \eta_{n-1}^{\beta-1} (\mathbf{u}_n - \mathbf{x}^*)$  converges to a non-zero vector

$$\frac{1}{\gamma} \left\{ \mathcal{P}_{A^\perp} \left( \nabla^2 f(\mathbf{x}^*) - \frac{1-\beta}{\eta_0} \mathbf{1}_{\{\alpha=1\}} \mathbf{I} \right) \mathcal{P}_{A^\perp} \right\}^\dagger \left( \mathcal{P}_{A^\perp} \nabla^2 f(\mathbf{x}^*) \nabla f(\mathbf{x}^*) \right)$$

- Remark: From the above derivation, for the choice  $p_n \propto \eta_n^\beta$ , when  $\beta$  varies, our algorithm has an interesting bias-variance tradeoff.



# Interesting Bias-Variance Tradeoff

- Order of fluctuation:  $\mathcal{O}(\eta_n^{1/2})$
- Order of Bias:  $\mathcal{O}(\eta_n^{1-\beta})$
- When  $\beta \in [0, 1/2)$   
The fluctuation caused by the randomness of gradient queries in every iteration dominates the optimization accuracy.
- When  $\beta \in [1/2, 1)$   
Manipulated by the biases formed by the accumulation of skewed updates in the unconstrained state within each 'inner loop'

**The End**