
Efficiency Ordering of Stochastic Gradient Descent – Supplementary Material

Anonymous Author(s)
Affiliation
Address
email

1	Contents	
2	A CLT for Stochastic Approximation Algorithms	2
3	B Proof of Lemma 3.1	3
4	C Discussion on Polyak-Lojasiewicz Inequality and Positive Definite Matrix $\nabla^2 f(\theta^*)$	5
5	D Proof of Theorem 3.6	5
6	D.1 Proof of Theorem 3.6 (i)	5
7	D.2 Proof of Theorem 3.6 (ii)	6
8	E Additional Convergence and CLT results for SGD with Constant Step Size and Quadratic Objective Function	6
9		
10	F Proof of Proposition 4.1	9
11	G Proof of Lemma 4.2	10
12	H Proof of Proposition 4.3 and 4.4	11
13	H.1 Single Shuffling in SGD CLT Analysis	11
14	H.2 Augmentation of Random Shuffling for CLT Analysis	12
15	H.3 Extension to Mini-batch Gradient Descent	13
16	I Simulation	15
17	I.1 Details behind Figure 1	15
18	I.2 Numerical Results on Large Graphs	19
19	I.3 Additional Simulations on Non-convex Objective Function and SGD Variants	19

20 In this supplementary material, any theorems, lemmas, and propositions that are indexed without
 21 citation are referred to our submission.

22 A CLT for Stochastic Approximation Algorithms

23 The existing central limit theorem (CLT) for stochastic approximation (SA) with Markovian dynamics
 24 [4, 11, 14] usually studied a general Markov process $\{X_t\}_{t \geq 0}$ on the finite state space \mathcal{V} and its
 25 transition kernel P_θ dependent on θ such that $P(X_{t+1} \in A | X_t = x, \theta_t = \theta) = P_\theta(x, A)$ for any
 26 subset $A \subseteq \mathcal{V}$. Denote π_θ the stationary distribution of P_θ . Define $P_\theta v_\theta(x) \triangleq \sum_{l \in \mathcal{V}} [P_\theta]_{x,l} \cdot v_\theta(l)$.
 27 The general SA algorithm is of the form

$$\theta_{t+1} = \text{Proj}_\Theta(\theta_t + \gamma_{t+1} H(\theta_t, X_{t+1})), \quad (15)$$

where $\Theta \subset \mathbb{R}^d$ is a closed and convex set. The main goal is to find the root θ^* of function

$$h(\theta) \triangleq \mathbb{E}_{X \sim \pi_\theta} [H(\theta, X)] \text{ i.e., } h(\theta^*) = 0.$$

28 As mentioned in [4] p.332 Theorem 13, and [11] p.31 Theorem 15, p.59 Theorem 25, the usual
 29 assumptions are given as

30 (B1) Function $h : \Theta \rightarrow \mathbb{R}^d$ is continuous on Θ , there exists a non-negative C^1 function V such
 31 that $\langle \nabla V(\theta), h(\theta) \rangle \leq 0, \forall \theta \in \Theta$ and the set $\mathcal{S} = \{\theta; \langle \nabla V(\theta), h(\theta) \rangle = 0\}$ is such that
 32 $V(\mathcal{S})$ has empty interior. Also, $V(\theta)$ tends to $+\infty$ if $\theta \rightarrow \partial\Theta$, where $\partial\Theta$ is the boundary of
 33 Θ , or $\|\theta\|_2 \rightarrow \infty$. There exists a compact set $\mathcal{K} \subset \Theta$ such that $\langle \nabla V(\theta), h(\theta) \rangle < 0$ if $\theta \notin \mathcal{K}$;

34 (B2) For every θ , there exist a function $v_\theta(x)$ such that the Poisson equation

$$v_\theta(x) - \mathbb{E}[v_\theta(X_{t+1}) | X_t = x, \theta_t = \theta] = H(\theta, X) - h(\theta). \quad (16)$$

35 For any compact set $\mathcal{C} \subset \Theta$,

$$\sup_{\theta \in \mathcal{C}, x \in \mathcal{V}} \|H(\theta, x)\|_2 + \|v_\theta(x)\|_2 < \infty. \quad (17)$$

36 There exists a continuous function $\phi_{\mathcal{C}}, \phi_{\mathcal{C}}(0) = 0$, such that for any $\theta, \theta' \in \mathcal{C}$,

$$\sup_{X \in \mathcal{V}} \|P_\theta v_\theta(x) - P_{\theta'} v_{\theta'}(x)\|_2 \leq \phi_{\mathcal{C}}(\|\theta - \theta'\|_2). \quad (18)$$

37 (B3) The step size follows $\gamma_t \geq 0, \sum_{t \geq 1} \gamma_t = \infty, \sum_{t \geq 1} \gamma_t^2 < \infty$ and $\sum_{t \geq 1} |\gamma_{t+1} - \gamma_t| < \infty$.

38 (B4) Assume θ_t converges to some limit $\theta^* \in \mathcal{S}$. Function h is C^1 in some neighborhood of θ^*
 39 with first derivatives Lipschitz, and matrix $\nabla h(\theta^*)$ has all its eigenvalues with negative real
 40 part.

41 Then, we have the following convergence and CLT result.

42 **Theorem A.1.** [4, 11, 14] Assume θ_t is given by the SA iteration (15) that satisfies assumptions (B1)
 43 – (B3) above, then iterate θ_t converges almost surely to the set \mathcal{S} defined in (B1). Moreover, with
 44 additional assumption (B4), we have

$$\frac{1}{\sqrt{\gamma_t}} \cdot (\theta_t - \theta^*) \xrightarrow[t \rightarrow \infty]{Dist} \mathcal{N}(0, \mathbf{V}_X), \quad (19)$$

45 where covariance matrix \mathbf{V}_X is the unique solution to the following Lyapunov equation:

$$\begin{cases} \Sigma_X + \mathbf{K} \mathbf{V}_X + \mathbf{V}_X \mathbf{K}^T = \mathbf{0} & \text{if } \alpha \in (\frac{1}{2}, 1), \\ \Sigma_X + (\mathbf{K} + \frac{\mathbf{I}}{2}) \mathbf{V}_X + \mathbf{V}_X (\mathbf{K} + \frac{\mathbf{I}}{2})^T = \mathbf{0} & \text{if } \alpha = 1. \end{cases} \quad (20)$$

46 Here, $\mathbf{K} \triangleq \nabla h(\theta^*)$ and $\Sigma_X \triangleq \Sigma_X(H(\theta^*, \cdot))$ is the asymptotic covariance matrix as in (8),
 47 evaluated at function $H(\theta^*, \cdot)$.

48 In addition, for averaged iterates $\bar{\theta}_t \triangleq \frac{1}{t} \sum_{i=0}^{t-1} \theta_i$, we still have $\bar{\theta}_t \xrightarrow[t \rightarrow \infty]{a.s.} \theta^*$, and

$$\sqrt{t} \cdot (\bar{\theta}_t - \theta^*) \xrightarrow[t \rightarrow \infty]{Dist} \mathcal{N}(0, \mathbf{V}'_X), \quad (21)$$

49 where $\mathbf{V}'_X = \mathbf{K}^{-1} \Sigma_X (\mathbf{K}^{-1})^T$ with the same matrices \mathbf{K} and Σ_X as in (20). \square

50 B Proof of Lemma 3.1

51 To prove Lemma 3.1 with existing Theorem A.1, we need to show that (A1) – (A5) is a special case
52 of (B1) – (B4). We list (A1) – (A5) here for self-contained purpose.

53 (A1) The step size is given by $\gamma_t = t^{-\alpha}$ for $\alpha \in (1/2, 1]$;

54 (A2) There exists a unique minimizer θ^* in the interior of the compact set Θ with $\nabla f(\theta^*) = 0$,
55 and matrix $\nabla^2 f(\theta^*)$ (resp. $\nabla^2 f(\theta^*) - \mathbf{I}/2$) is positive definite for $a \in (1/2, 1)$ (resp. $a = 1$);

56 (A3) Gradients are bounded in the compact set Θ , that is, $\sup_{\theta \in \Theta} \sup_{i \in [n]} \|\nabla F(\theta, i)\|_2 < \infty$;

57 (A4) For every $z \in [n], \theta \in \mathbb{R}^d$, the solution $\tilde{F}(\theta, z) \in \mathbb{R}^d$ of the Poisson equation

$$\tilde{F}(\theta, z) - \mathbb{E}[\tilde{F}(\theta, X_{t+1}) \mid X_t = z] = \nabla F(\theta, z) - \nabla f(\theta) \quad (22)$$

58 exists, and $\sup_{\theta \in \Theta, z \in [n]} \|\tilde{F}(\theta, z)\|_2 < \infty$;

59 (A5) The functions $F(\theta, i)$ are L -smooth for all $i \in [n]$, that is, $\forall \theta_1, \theta_2 \in \Theta, \forall i \in [n]$, we have
60 $\|\nabla F(\theta_1, i) - \nabla F(\theta_2, i)\|_2 \leq L\|\theta_1 - \theta_2\|_2$.

61 Let $H(\theta, X) \triangleq -\nabla F(\theta, X)$ for function $F(\theta, X)$ defined in (1). Then, we have $h(\theta) \triangleq$
62 $\mathbb{E}_{X \sim \pi}[H(\theta, X)] = -\nabla f(\theta)$. By choosing $V(\theta) \triangleq f(\theta)$, we know $\langle \nabla V(\theta), h(\theta) \rangle = -\nabla f(\theta)^2 \leq 0$.
63 From (A2) we know θ^* is the unique minimizer of function f , by letting $\mathcal{K} = \{\theta^*\}$, we have
64 $\langle \nabla V(\theta), h(\theta) \rangle < 0$ when $\theta \notin \mathcal{K}$. Therefore, **(B1) is satisfied**.

65 Now we need to check assumption (B2). Assumption (A4) is a direct translation to (16) in (B2),
66 and $\sup_{\theta \in \Theta, z \in [n]} \|\tilde{F}(\theta, z)\|_2 < \infty$, as well as assumption (A3), implies (17). We still need to show
67 (18). By assuming an n -state ergodic Markov chain $\{X_t\}_{t \geq 0}$ (θ -independent) with transition kernel
68 $\mathbf{P} \in \mathbb{R}^{n \times n}$ and stationary distribution π , the solution $\tilde{F}(\theta, z)$ to the Poisson equation (22) in (A4)
69 exists and is given as follows.¹

$$\tilde{F}(\theta, z) = \nabla F(\theta, z) - \nabla f(\theta) + \sum_{l=1}^n \mathbf{P}_{z,l} (\nabla F(\theta, l) - \nabla f(\theta)) + \sum_{l=1}^n [\mathbf{P}^2]_{z,l} (\nabla F(\theta, l) - \nabla f(\theta)) + \dots \quad (23)$$

70 Next, we can rewrite $\tilde{F}(\theta, z)$ in the closed form and show that it is Lipschitz continuous and satisfies
71 (20) in assumption (B2). Note that by definition of expectation and Chapman–Kolmogorov equation
72 ($\sum_{k=1}^n \sum_{l=1}^n \mathbf{P}_{z,k} \mathbf{P}_{k,l} = \sum_{l=1}^n [\mathbf{P}^2]_{z,l}$), we have

$$\begin{aligned} \mathbb{E}[\tilde{F}(\theta, X_{t+1}) \mid X_t = z] &= \sum_{l=1}^n \mathbf{P}_{z,l} (\nabla F(\theta, l) - \nabla f(\theta)) + \sum_{l=1}^n [\mathbf{P}^2]_{z,l} (\nabla F(\theta, l) - \nabla f(\theta)) \\ &\quad + \sum_{l=1}^n [\mathbf{P}^3]_{z,l} (\nabla F(\theta, l) - \nabla f(\theta)) + \dots \end{aligned} \quad (24)$$

73 Then, from (23) and (24) we have $\tilde{F}(\theta, z) - \mathbb{E}[\tilde{F}(\theta, X_{t+1}) \mid X_t = z] = \nabla F(\theta, z) - \nabla f(\theta)$, which
74 is exactly (22) in (A4). Moreover, since $\mathbf{1}$ and π are the right and left eigenvectors of \mathbf{P} respectively
75 with eigenvalue 1, by induction we know

$$\mathbf{P}^k - \mathbf{1}\pi^T = (\mathbf{P} - \mathbf{1}\pi^T)^k, \forall k \in \mathbb{Z}, k \geq 1. \quad (25)$$

76 Along with the fact that

$$\nabla f(\theta) = \sum_{l=1}^n \pi_l \nabla F(\theta, l) = \sum_{l=1}^n [\mathbf{1}\pi^T]_{z,l} \nabla F(\theta, l), \quad (26)$$

¹In this paper, we only consider θ -independent Markovian inputs. The more general conditions of the θ -dependent Markov chain under which the solution of (22) in (A4) exists are referred to [11] p.71 Theorem 35 or [4] p.217.

77 we can further simplify and get a closed form of (23), which is given below.

$$\begin{aligned}
\tilde{F}(\theta, z) &= \nabla F(\theta, z) - \nabla f(\theta) + \sum_{l=1}^n \mathbf{P}_{z,l} (\nabla F(\theta, l) - \nabla f(\theta)) + \sum_{l=1}^n [\mathbf{P}^2]_{z,l} (\nabla F(\theta, l) - \nabla f(\theta)) + \dots \\
&= \sum_{l=1}^n [\mathbf{P}^0 - \mathbf{1}\pi^T]_{z,l} \nabla F(\theta, l) + \sum_{l=1}^n [\mathbf{P}^1 - \mathbf{1}\pi^T]_{z,l} \nabla F(\theta, l) + \dots \\
&= \left(\sum_{k=0}^{\infty} \sum_{l=1}^n [\mathbf{P}^k - \mathbf{1}\pi^T]_{z,l} \nabla F(\theta, l) \right) - \nabla f(\theta) \\
&= \left(\sum_{l=1}^n \left[\sum_{k=0}^{\infty} [\mathbf{P}^k - \mathbf{1}\pi^T]_{z,l} \right] \nabla F(\theta, l) \right) - \nabla f(\theta) \\
&= \left(\sum_{l=1}^n \left[\sum_{k=0}^{\infty} [(\mathbf{P} - \mathbf{1}\pi^T)^k]_{z,l} \right] \nabla F(\theta, l) \right) - \nabla f(\theta) \\
&= \sum_{l=1}^n \left[(\mathbf{I} - \mathbf{P} + \mathbf{1}\pi^T)^{-1} \right]_{z,l} \nabla F(\theta, l) - \nabla f(\theta),
\end{aligned} \tag{27}$$

78 where the second equality comes from (26) and the fifth equality is from (25). Recall the definition

79 $\mathbf{P}F(\theta, z) \triangleq \sum_{l=1}^n \mathbf{P}_{z,l} F(\theta, l)$, we can show that

$$\begin{aligned}
&\|\mathbf{P}\tilde{F}(\theta, z) - \mathbf{P}\tilde{F}(\theta', z)\|_2 \\
&= \left\| \sum_{l=1}^n \mathbf{P}_{z,l} (\tilde{F}(\theta, l) - \tilde{F}(\theta', l)) \right\|_2 \\
&\leq \sum_{l=1}^n \mathbf{P}_{z,l} \|\tilde{F}(\theta, l) - \tilde{F}(\theta', l)\|_2 \\
&\leq \sup_{z \in \mathcal{V}} \|\tilde{F}(\theta, z) - \tilde{F}(\theta', z)\|_2 \\
&\leq \sup_{z \in \mathcal{V}} \left\| \sum_{l=1}^n [(\mathbf{I} - \mathbf{P} + \mathbf{1}\pi^T)^{-1}]_{z,l} (\nabla F(\theta, l) - \nabla F(\theta', l)) \right\|_2 + \|\nabla f(\theta) - \nabla f(\theta')\|_2 \\
&\leq C \sup_{z \in \mathcal{V}} \|\nabla F(\theta, z) - \nabla F(\theta', z)\|_2 + \|\nabla f(\theta) - \nabla f(\theta')\|_2 \\
&\leq (C+1)L\|\theta - \theta'\|_2
\end{aligned} \tag{28}$$

for some constant C related to matrix $(\mathbf{I} - \mathbf{P} + \mathbf{1}\pi^T)^{-1}$, where the first and the third inequalities are from triangular inequality and the last inequality comes from assumption (A5). Note that we have

$$\|\nabla f(\theta) - \nabla f(\theta')\|_2 = \left\| \sum_{i=1}^n \pi_i (\nabla F(\theta, i) - \nabla F(\theta', i)) \right\|_2 \leq \sup_{z \in \mathcal{V}} \|\nabla F(\theta, i) - \nabla F(\theta', i)\|_2 \leq L\|\theta - \theta'\|_2$$

80 in the last inequality of (28). So (18) is shown and (B2) is satisfied.

81 For assumption (A1) with respect to the conditions on the step size, we know for $a \in (1/2, 1]$,
82 $\sum_{t \geq 1} 1/t^a = \infty$ and $\sum_{t \geq 0} 1/t^{2a} < \infty$. Besides,

$$\gamma_t - \gamma_{t+1} = \frac{1}{t^a} - \frac{1}{(t+1)^a} = \frac{(t+1)^a - t^a}{t^a(t+1)^a} \leq \frac{(t+1)^a - t^a}{t^{2a}} \leq \frac{1}{t^{2a}}$$

83 where the second inequality comes from $(t+1)^a - t^a$ monotone decreasing in t for $a \in (1/2, 1]$.
84 Then, we have $\sum_{t \geq 1} |\gamma_t - \gamma_{t+1}| \leq \sum_{t \geq 1} 1/t^{2a} < \infty$. Then, (B3) is satisfied.

85 Since (B1) – (B3) are satisfied and (A2) assumes unique minimizer such that $\mathcal{K} = \{\theta^*\}$, from
86 Theorem A.1 we know $\theta_t \xrightarrow[t \rightarrow \infty]{a.s.} \theta^*$. Along with assumption (A2) on the positive definite matrix
87 $\nabla^2 f(\theta^*)$, (B4) is satisfied.

88 Therefore, (A1)-(A5) implies (B1) – (B4) and all the results from Theorem A.1 can be carried over to
 89 Lemma 3.1.

90 C Discussion on Polyak-Lojasiewicz Inequality and Positive Definite Matrix 91 $\nabla^2 f(\theta^*)$

92 In this part, we discuss the strictness of the condition on the objective function f between Polyak-
 93 Lojasiewicz (P-L) inequality and our assumption (A2) - positive definite matrix $\nabla^2 f(\theta^*)$. We say that
 94 if a scalar-valued function f satisfies μ -P-L inequality, then for any $\theta \in \mathbb{R}^d$, the following condition
 95 holds:

$$\frac{1}{2} \|\nabla f(\theta)\|_2^2 \geq \mu(f(\theta) - f(\theta^*)), \quad (29)$$

where $\nabla f(\theta) \in \mathbb{R}^d$, $f(\theta^*) = \min_{\theta \in \mathbb{R}^d} f(\theta)$ and the minimizer θ^* belongs to a non-empty solution
 set. We define a new function

$$g(\theta) \triangleq \frac{1}{2} \|\nabla f(\theta)\|_2^2 - \mu(f(\theta) - f^*).$$

96 Then, (29) is equivalent to saying $\min_{\theta} g(\theta) \geq 0$, and the necessary condition to ensure that θ^* is the
 97 local minimizer is $\nabla^2 g(\theta^*) \succeq_L \mathbf{0}$ (e.g., Chapter 1.2 [25]). We have

$$\nabla^2 g(\theta) = (\nabla^2 f(\theta) - \mu \mathbf{I}) \nabla^2 f(\theta) + \mathbf{M} \otimes \nabla f(\theta),$$

98 where matrix \mathbf{M} is a 3D matrix with dimension $d \times d \times d$ and \otimes is the tensor product. Since $\nabla f(\theta^*) =$
 99 0 , we have $\mathbf{M} \otimes (\nabla f(\theta^*)) = 0$. Then, $\nabla^2 g(\theta^*) \succeq_L \mathbf{0}$ implies $(\nabla^2 f(\theta^*))^2 \succeq_L \mu \nabla^2 f(\theta^*)$. Denote
 100 $\lambda_i \geq 0, i = 1, 2, \dots, d$ the eigenvalues of matrix $\nabla^2 f(\theta^*)$, by spectral decomposition we need
 101 $\lambda_i \geq \mu$ or $\lambda_i = 0$ for each i . If all the eigenvalues of $\nabla^2 f(\theta^*)$ are no smaller than μ , then $\nabla^2 f(\theta^*)$ is
 102 a positive definite matrix by definition. For example, μ -strongly convex objective function f satisfies
 103 both P-L inequality and $\nabla^2 f(\theta^*)$ being positive definite. If there exists at least one eigenvalue with
 104 zero value, then $\nabla^2 f(\theta^*)$ is no longer positive definite.

On the other hand, positive definite matrix $\nabla^2 f(\theta^*)$ does not necessarily imply P-L inequality. We
 give a toy example of objective function f that satisfies positive definite matrix $\nabla^2 f(\theta^*)$ while fails
 to satisfy P-L inequality. For some smooth convex function

$$f(\theta) = \sqrt{\|\theta\|_2^2 + 1} \geq 1,$$

we know

$$f'(\theta) = \frac{\theta}{\sqrt{\|\theta\|_2^2 + 1}}, \quad f''(\theta) = \frac{1}{\sqrt{\|\theta\|_2^2 + 1}} \mathbf{I} - \frac{1}{(\|\theta\|_2^2 + 1)^{3/2}} \theta \theta^T$$

for any $\theta \in \mathbb{R}^d$. Since $\theta^* = \mathbf{0}$, $f(\theta^*) = 1$ and $f''(\theta^*) = \mathbf{I}$ is a positive definite matrix such that this
 objective function satisfies our assumption (A2). However, for any $\theta \in \mathbb{R}^d$, there always exists a
 constant $\epsilon > 0$ such that

$$\epsilon(f(\theta) - f(\theta^*)) \geq \|f'(\theta)\|_2^2,$$

105 which fails to satisfy (29). Therefore, there is no inclusive relationship between P-L inequality and
 106 positive definite matrix $\nabla^2 f(\theta^*)$. Both of the conditions can cover different types of functions.

107 D Proof of Theorem 3.6

108 D.1 Proof of Theorem 3.6 (i)

109 We first prove the direction that efficiency ordering implies Loewner ordering. For any vector
 110 $\mathbf{v} \triangleq [v_1, v_2, \dots, v_d]^T \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ and vector-valued function $\mathbf{f}(X) \triangleq [f_1(X), f_2(X), \dots, f_d(X)]^T$,

111 with $\Sigma(\mathbf{f})$ defined in (8) we can get

$$\begin{aligned}
\mathbf{v}^T \Sigma(\mathbf{f}) \mathbf{v} &= \lim_{t \rightarrow \infty} \mathbf{v}^T \Sigma(\mathbf{f}, t) \mathbf{v} \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left\{ \mathbf{v}^T \left[\sum_{s=1}^t (\mathbf{f}(X_s) - \mathbb{E}_\pi[\mathbf{f}(X)]) \right] \left[\sum_{s=1}^t (\mathbf{f}(X_s) - \mathbb{E}_\pi[\mathbf{f}(X)]) \right]^T \mathbf{v} \right\} \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left\{ \left[\sum_{s=1}^t (g_{\mathbf{v}, \mathbf{f}}(X_s) - \mathbb{E}_\pi[g_{\mathbf{v}, \mathbf{f}}(X)]) \right]^2 \right\} \\
&= \sigma^2(g_{\mathbf{v}, \mathbf{f}}),
\end{aligned} \tag{30}$$

where function

$$g_{\mathbf{v}, \mathbf{f}}(X) \triangleq v_1 f_1(X) + v_2 f_2(X) + \dots + v_d f_d(X)$$

112 is a linear combination of $f_i(X)$. For two random processes with efficiency ordering and an arbitrary
113 vector-valued function \mathbf{f} , $\sigma_X^2(g_{\mathbf{v}, \mathbf{f}}) \leq \sigma_Y^2(g_{\mathbf{v}, \mathbf{f}})$ for any vector \mathbf{v} , which is exactly $\mathbf{v}^T \Sigma_X(\mathbf{f}) \mathbf{v} \leq$
114 $\mathbf{v}^T \Sigma_Y(\mathbf{f}) \mathbf{v}$. Then, by definition of Loewner ordering, we have $\Sigma_X(\mathbf{f}) \leq_L \Sigma_Y(\mathbf{f})$ for any vector-
115 valued function \mathbf{f} .

116 On the other direction, let $\mathbf{v} = [1, 0, \dots, 0]^T$ and vector-valued function $\mathbf{f}(X) = [g(X), 0, \dots, 0]^T$,
117 where g can be any scalar-valued function. Then, (30) can be written as $\mathbf{v}^T \Sigma(\mathbf{f}) \mathbf{v} = \sigma^2(g)$ and it
118 holds for any scalar-valued function g . For two Markov chains $\{X_t\}, \{Y_t\}$ with $\Sigma_X(\mathbf{f}) \leq_L \Sigma_Y(\mathbf{f})$
119 for any vector-valued function \mathbf{f} , we have $\mathbf{v}^T \Sigma_X(\mathbf{f}) \mathbf{v} \leq \mathbf{v}^T \Sigma_Y(\mathbf{f}) \mathbf{v}$ for any vector \mathbf{v} . Then, with
120 $\mathbf{v} = [1, 0, \dots, 0]^T$ we show that $\sigma_X^2(g) \leq \sigma_Y^2(g)$ for any scalar-valued function g , which proves the
121 efficiency ordering.

122 D.2 Proof of Theorem 3.6 (ii)

123 We first introduce the closed form of the solution \mathbf{V} to the Lyapunov equation in Lemma 3.1 and the
124 useful lemma on Loewner ordering.

125 **Lemma D.1** ([8] Theorem 3.16 and (3.160)). *If all the eigenvalues of matrix \mathbf{K} have negative*
126 *real part, then for every positive-definite matrix \mathbf{U} there exists a unique positive-definite matrix \mathbf{V}*
127 *satisfying $\mathbf{U} + \mathbf{K}\mathbf{V} + \mathbf{V}\mathbf{K}^T = \mathbf{0}$. The explicit solution \mathbf{V} is given as*

$$\mathbf{V} = \int_0^\infty e^{\mathbf{K}t} \mathbf{U} e^{(\mathbf{K}^T)t} dt. \tag{31}$$

128 **Lemma D.2** ([28] Theorem 8.2.7). *If two real matrix $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times m}$ are Loewner ordered $\mathbf{A} \leq_L \mathbf{B}$,*
129 *then $\mathbf{C}\mathbf{A}\mathbf{C}^T \leq_L \mathbf{C}\mathbf{B}\mathbf{C}^T$ for any real matrix $\mathbf{C} \in \mathbb{R}^{m \times m}$.*

130 From Theorem 3.6 (i), we know efficiency ordering $\sigma_X^2(g) \leq \sigma_Y^2(g)$ for any scalar-valued function
131 g leads to Loewner ordering $\Sigma_X(\mathbf{f}) \leq_L \Sigma_Y(\mathbf{f})$ for any vector-valued function \mathbf{f} . Consider two
132 random process $\{X_t\}_{t \geq 0}, \{Y_t\}_{t \geq 0}$ with efficiency ordering $X \geq_E Y$, we have $\Sigma_X \leq_L \Sigma_Y$. By
133 Lemma D.2 and (20), for any t in (31), we have $e^{\mathbf{K}t} \Sigma_X e^{(\mathbf{K}^T)t} \leq_L e^{\mathbf{K}t} \Sigma_Y e^{(\mathbf{K}^T)t}$. Then, for any
134 vector $\mathbf{v} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$, we have

$$\mathbf{v}^T \mathbf{V}_X \mathbf{v} = \int_0^\infty \mathbf{v}^T e^{\mathbf{K}t} \Sigma_X e^{(\mathbf{K}^T)t} \mathbf{v} dt \leq \int_0^\infty \mathbf{v}^T e^{\mathbf{K}t} \Sigma_Y e^{(\mathbf{K}^T)t} \mathbf{v} dt = \mathbf{v}^T \mathbf{V}_Y \mathbf{v},$$

135 such that $\mathbf{V}_X \leq_L \mathbf{V}_Y$ by definition of Loewner ordering. Similarly, for averaged iterates, we have
136 $\mathbf{V}'_X \leq_L \mathbf{V}'_Y$ immediately from Lemma D.2 because $\Sigma_X \leq_L \Sigma_Y$ and $\mathbf{V}'_X = \mathbf{K}^{-1} \Sigma_X (\mathbf{K}^{-1})^T$,
137 $\mathbf{V}'_Y = \mathbf{K}^{-1} \Sigma_Y (\mathbf{K}^{-1})^T$.

138 E Additional Convergence and CLT results for SGD with Constant Step Size 139 and Quadratic Objective Function

140 Lemma 3.1 has shown the CLT result for general SGD iteration (7) with diminishing step size. A
141 natural question would be if any CLT result exists for the same SGD iteration with constant step size
142 γ . For the *i.i.d* inputs and a special case of the iteration

$$\theta_{t+1} = \theta_t - \gamma(\mathbf{A}(X_{t+1})\theta_t - b(X_{t+1})), \tag{32}$$

143 which is usually called linear stochastic approximation in the stochastic approximation literature, it
 144 has been studied in [13, 29] that θ_t forms a Markov chain and its time-averaged iterate $\bar{\theta}_t = \frac{1}{t} \sum_{i=0}^{t-1} \theta_i$
 145 converges to the minimizer θ^* almost surely and a CLT result is given in Theorem 1 [29]. However,
 146 for Markovian inputs $\{X_t\}_{t \geq 0}$, V. Borkar and S. Meyn mentioned in [5] that the behavior of (θ_t, X_t)
 147 itself is still an open problem under the SGD iteration (7). In this part, we propose Lemma E.1
 148 that studies the special case of the SGD iteration [9] (which studied the diminishing step size) and
 149 complements the CLT result for constant step size with time-averaged iterates.

We consider a quadratic objective function

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{2} \theta^T \mathbf{A} \theta - \theta^T \mathbf{b}(i) \right),$$

150 where matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is positive definite and vector $\mathbf{b}(X) \in \mathbb{R}^d$ only depends on the state $X \in \mathcal{V}$
 151 of the Markovian input. Then, the SGD iteration studied in [9] is given as

$$\theta_{t+1} = \theta_t - \gamma_{t+1} (\mathbf{A} \theta_t - \mathbf{b}(X_{t+1})). \quad (33)$$

152 Here, we study the constant step size $\gamma_t = \gamma$, $\forall t \geq 0$. Define $\bar{\mathbf{b}} \triangleq \sum_{i \in [n]} \mathbf{b}(X_i) \pi_i$. The minimizer
 153 is given by $\theta^* = \mathbf{A}^{-1} \bar{\mathbf{b}}$ such that $\nabla f(\theta^*) = 0$. Then, we have the following CLT result for the SGD
 154 update rule (33) with constant step size and Markovian input $\{X_t\}_{t \geq 0}$.

155 **Lemma E.1.** Consider the update rule (33) with positive definite matrix \mathbf{A} and constant step size γ
 156 such that $0 < \gamma < 2/\|\mathbf{A}\|_2$. Then, for averaged iterates $\bar{\theta}_t = \frac{1}{t} \sum_{i=0}^{t-1} \theta_i$, we have

$$\bar{\theta}_t \xrightarrow[t \rightarrow \infty]{a.s.} \theta^*, \quad \text{and} \quad \sqrt{t}(\bar{\theta}_t - \theta^*) \xrightarrow[t \rightarrow \infty]{Dist} \mathcal{N}(0, \mathbf{V}_X), \quad (34)$$

157 where $\mathbf{V}_X = \mathbf{A}^{-1} \Sigma_X (\mathbf{A}^{-1})^T$ and $\Sigma_X = \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E}[B_t B_t^T]$, $B_t \triangleq \sum_{s=1}^t (\mathbf{b}(X_s) - \bar{\mathbf{b}})$.

158 *Proof.* Let $\tilde{\theta}_t = \theta_t - \theta^*$ and recall $\theta^* = \mathbf{A}^{-1} \bar{\mathbf{b}}$, we can rewrite (33) as

$$\tilde{\theta}_{t+1} = \tilde{\theta}_t - \gamma (\mathbf{A} \tilde{\theta}_t - \mathbf{b}(X_{t+1}) + \bar{\mathbf{b}}). \quad (35)$$

159 Recursively solving (35) gives

$$\tilde{\theta}_t = (\mathbf{I} - \gamma \mathbf{A})^t \tilde{\theta}_0 - \gamma \sum_{i=1}^t (\mathbf{I} - \gamma \mathbf{A})^{t-i} (\mathbf{b}(X_i) - \bar{\mathbf{b}}). \quad (36)$$

160 For averaged iterates $\bar{\theta}_t = \frac{1}{t} \sum_{i=0}^{t-1} \theta_i$, (36) gives

$$\begin{aligned} \bar{\theta}_t - \theta^* &= \frac{1}{t} \sum_{i=0}^{t-1} \tilde{\theta}_i \\ &= \frac{1}{t} \sum_{i=0}^{t-1} (\mathbf{I} - \gamma \mathbf{A})^i \tilde{\theta}_0 - \frac{\gamma}{t} \sum_{i=1}^{t-1} \sum_{j=1}^i (\mathbf{I} - \gamma \mathbf{A})^{i-j} (\mathbf{b}(X_j) - \bar{\mathbf{b}}) \\ &= \frac{1}{t} \sum_{i=0}^{t-1} (\mathbf{I} - \gamma \mathbf{A})^i \tilde{\theta}_0 - \frac{\gamma}{t} \sum_{i=1}^{t-1} \left[\sum_{j=0}^{t-i-1} (\mathbf{I} - \gamma \mathbf{A})^j \right] (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \\ &= \frac{1}{t} (\gamma \mathbf{A})^{-1} (\mathbf{I} - (\mathbf{I} - \gamma \mathbf{A})^t) \tilde{\theta}_0 - \frac{\gamma}{t} \sum_{i=1}^{t-1} (\gamma \mathbf{A})^{-1} (\mathbf{I} - (\mathbf{I} - \gamma \mathbf{A})^{t-i}) (\mathbf{b}(X_i) - \bar{\mathbf{b}}), \end{aligned} \quad (37)$$

161 where the third equality comes from rearranging the summation order in the second term on the RHS.
 162 The fourth equality comes from the fact that $\sum_{i=0}^{t-1} (\mathbf{I} - \gamma \mathbf{A})^i = (\gamma \mathbf{A})^{-1} (\mathbf{I} - (\mathbf{I} - \gamma \mathbf{A})^t)$.

163 Next we want to show $\lim_{t \rightarrow \infty} (\mathbf{I} - \gamma \mathbf{A})^t = \mathbf{0}$. Since we assume $0 < \gamma < 2/\|\mathbf{A}\|_2$, we have
 164 $\|\mathbf{I} - \gamma \mathbf{A}\|_2 = \max_{i=1,2,\dots,n} |1 - \gamma \lambda_i(\mathbf{A})| < 1$, where $\lambda_i(\mathbf{A}) > 0$ is the i -th eigenvalue of the
 165 positive definite matrix \mathbf{A} . Then, by submultiplicative property, $\|(\mathbf{I} - \gamma \mathbf{A})^t\|_2 \leq \|\mathbf{I} - \gamma \mathbf{A}\|_2^t$ such
 166 that $\lim_{t \rightarrow \infty} \|(\mathbf{I} - \gamma \mathbf{A})^t\|_2 \leq \lim_{t \rightarrow \infty} \|\mathbf{I} - \gamma \mathbf{A}\|_2^t = 0$, which implies that $\lim_{t \rightarrow \infty} (\mathbf{I} - \gamma \mathbf{A})^t = \mathbf{0}$.

167 Now we want to show $\lim_{t \rightarrow \infty} \left\| \sum_{i=1}^{t-1} (\mathbf{I} - \gamma \mathbf{A})^{t-i} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \right\|_2 < \infty$. Since vector-valued
 168 function $\mathbf{b}(\cdot)$ is defined on the finite state space \mathcal{V} , it is safe to assume $\|\mathbf{b}(X) - \bar{\mathbf{b}}\|_2 \leq C$ for some
 169 constant C . Then,

$$\begin{aligned} \lim_{t \rightarrow \infty} \left\| \sum_{i=1}^{t-1} (\mathbf{I} - \gamma \mathbf{A})^{t-i} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \right\|_2 &\leq \lim_{t \rightarrow \infty} \sum_{i=1}^{t-1} \|(\mathbf{I} - \gamma \mathbf{A})^{t-i}\|_2 \|\mathbf{b}(X_i) - \bar{\mathbf{b}}\|_2 \\ &\leq C \lim_{t \rightarrow \infty} \sum_{i=1}^{t-1} \|\mathbf{I} - \gamma \mathbf{A}\|_2^{t-i} \\ &= C \lim_{t \rightarrow \infty} \sum_{i=1}^{t-1} \|\mathbf{I} - \gamma \mathbf{A}\|_2^i < \infty, \end{aligned} \quad (38)$$

170 where the first inequality comes from submultiplicative property and triangular inequality, the first
 171 equality is by rewriting the index inside the summation, and the third inequality comes from the fact
 172 that $\|\mathbf{I} - \gamma \mathbf{A}\|_2 < 1$. Then, we have

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{i=1}^{t-1} (\mathbf{I} - \gamma \mathbf{A})^{t-i} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) = \mathbf{0}, \quad (39)$$

173 and

$$\lim_{t \rightarrow \infty} \frac{1}{\sqrt{t}} \sum_{i=1}^{t-1} (\mathbf{I} - \gamma \mathbf{A})^{t-i} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) = \mathbf{0}. \quad (40)$$

174 With $\lim_{t \rightarrow \infty} (\mathbf{I} - \gamma \mathbf{A})^t = \mathbf{0}$ and (39), we have from (37) that

$$\lim_{t \rightarrow \infty} \bar{\theta}_t - \theta^* = \lim_{t \rightarrow \infty} -\frac{\gamma}{t} \sum_{i=1}^{t-1} (\gamma \mathbf{A})^{-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) = -\mathbf{A}^{-1} \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{i=1}^{t-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}). \quad (41)$$

175 From the ergodic theorem for Markov chains ([7] Theorem 3.3.2), we have $\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{i=1}^{t-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) = \mathbf{0}$
 176 and therefore $\lim_{t \rightarrow \infty} \bar{\theta}_t = \theta^*$.

177 To get the CLT result in (34), we first scale $\bar{\theta}_t - \theta^*$ from (37), along with (40), such that

$$\begin{aligned} \lim_{t \rightarrow \infty} \sqrt{t} (\bar{\theta}_t - \theta^*) &= \lim_{t \rightarrow \infty} -\frac{\gamma}{\sqrt{t}} \sum_{i=1}^{t-1} (\gamma \mathbf{A})^{-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \\ &= -\mathbf{A}^{-1} \lim_{t \rightarrow \infty} \frac{\sqrt{t-1}}{\sqrt{t}} \left(\frac{1}{\sqrt{t-1}} \sum_{i=1}^{t-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \right). \end{aligned} \quad (42)$$

178 From the CLT of Markov chain in Theorem 2.1, we know $\frac{1}{\sqrt{t}} \sum_{i=1}^t (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \xrightarrow[t \rightarrow \infty]{dist} \mathcal{N}(0, \boldsymbol{\Sigma}_X)$,
 179 where $\boldsymbol{\Sigma}_X = \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E}[(\sum_{s=1}^t (\mathbf{b}(X_s) - \bar{\mathbf{b}})) (\sum_{s=1}^t (\mathbf{b}(X_s) - \bar{\mathbf{b}}))^T]$. This result shows that time-
 180 averaged iterate $\bar{\theta}_t$ will guarantee the convergence to the exact solution and we have CLT result for
 181 $\sqrt{t}(\bar{\theta}_t - \theta^*)$ too.

182 Finally, we need to quantify the covariance matrix in the CLT result to $\sqrt{t}(\bar{\theta}_t - \theta^*)$. We will look
 183 at $\lim_{t \rightarrow \infty} t \mathbb{E}[(\bar{\theta}_t - \theta^*) (\bar{\theta}_t - \theta^*)^T]$. Note that the second term in (37) is bounded (see (38) for the
 184 proof) such that the cross term in the outer-product of $\bar{\theta}_t - \theta^*$ will vanish when $t \rightarrow \infty$. Then, we
 185 have

$$\begin{aligned} &\lim_{t \rightarrow \infty} t (\bar{\theta}_t - \theta^*) (\bar{\theta}_t - \theta^*)^T \\ &= \lim_{t \rightarrow \infty} \frac{t-1}{t} \left(-\frac{1}{\sqrt{t-1}} \mathbf{A}^{-1} \sum_{i=1}^{t-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \right) \left(-\frac{1}{\sqrt{t-1}} \mathbf{A}^{-1} \sum_{i=1}^{t-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \right)^T \\ &= \mathbf{A}^{-1} \lim_{t \rightarrow \infty} \frac{1}{t} \left(\sum_{i=1}^{t-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \right) \left(\sum_{i=1}^{t-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \right)^T (\mathbf{A}^{-1})^T. \end{aligned} \quad (43)$$

186 Taking the expectation of (43) gives

$$\begin{aligned}
& \lim_{t \rightarrow \infty} t \mathbb{E}[(\bar{\theta}_t - \theta^*)(\bar{\theta}_t - \theta^*)^T] \\
&= \mathbf{A}^{-1} \lim_{t \rightarrow \infty} \frac{t-1}{t} \mathbb{E} \left[\frac{1}{t-1} \left(\sum_{i=1}^{t-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \right) \left(\sum_{i=1}^{t-1} (\mathbf{b}(X_i) - \bar{\mathbf{b}}) \right)^T \right] (\mathbf{A}^{-1})^T \quad (44) \\
&= \mathbf{A}^{-1} \Sigma_X (\mathbf{A}^{-1})^T.
\end{aligned}$$

Therefore, we have

$$\sqrt{t}(\bar{\theta}_t - \theta^*) \xrightarrow[t \rightarrow \infty]{dist} \mathcal{N}(0, \mathbf{A}^{-1} \Sigma_X (\mathbf{A}^{-1})^T).$$

187

□

188 In Lemma E.1, $\mathbf{A} = \nabla^2 f(\theta)$ and Σ_X , by definition (8), is an asymptotic covariance matrix of the
189 Markov chain $\{X_t\}_{t \geq 0}$ for vector-valued function $\mathbf{b}(\cdot)$. Therefore, (34) shares a similar form to (21)
190 in Lemma 3.1. Our Theorem 3.6 can be carried over to Lemma E.1, which enables us to compare the
191 efficiency ordering of SGD algorithms driven by different stochastic inputs under the update rule (33)
192 and constant step size.

193 F Proof of Proposition 4.1

[30, 18] proposed the guidance by modifying a reversible random walk into a non-Markovian random walk to achieve higher sampling efficiency and it was applied to other applications to improve sampling efficiency (e.g., [22, 24]). Specifically speaking, consider a reversible random walk $\{X_t\}_{t \geq 0}$ (e.g., SRW) with transition matrix \mathbf{P} and stationary distribution π . Let its counterpart (e.g., NBRW) on the augmented state space be given by $\{Z_t\}_{t \geq 0} \triangleq \{(Y_{t-1}, Y_t)\}_{t \geq 0}$, where $Y_{t-1}, Y_t \in \mathcal{V}$ and $Z_0 = (Y_0, Y_0)$. Additionally, $\{Z_t\}_{t \geq 0}$ is a Markov chain on the augmented state space

$$\mathcal{E} \triangleq \{(i, j) : i, j \in \mathcal{V} \text{ s.t. } P(i, j) > 0\} \subseteq \mathcal{V} \times \mathcal{V}$$

194 with stationary distribution π' . For notation simplicity, we use e_{ij} to represent edge (i, j) . Note that
195 by definition $e_{ij} \neq e_{ji}$ and we allow $i = j$ if $P(i, j) > 0$, which is a bit different from the edge set
196 that does not include edge (i, i) . As proved in [18], the properties of NBRW $\{Z_t\}_{t \geq 0}$ are detailed in
197 the following theorem.

198 **Theorem F.1** ([30] Theorem 2). *Suppose that $\{X_t\}$ is an irreducible, reversible Markov chain on
199 the state space $\mathcal{V} = \{1, 2, \dots, n\}$ with transition matrix $\mathbf{P} = \{P(i, j)\}$ and stationary distribution
200 π . Construct a Markov chain $\{Z_t\}$ on the augmented state space \mathcal{E} with transition matrix $\mathbf{P}' =$
201 $\{P'(e_{ij}, e_{lk})\}$ in which the transition probabilities $P'(e_{ij}, e_{lk})$ satisfy the following two conditions:
202 for all $e_{ij}, e_{ji}, e_{jk}, e_{kj} \in \mathcal{E}$ with $i \neq k$,*

$$P(j, i)P'(e_{ij}, e_{jk}) = P(j, k)P'(e_{kj}, e_{ji}), \quad (45a)$$

203

$$P'(e_{ij}, e_{jk}) \geq P(j, k). \quad (45b)$$

204 Then, the Markov chain $\{Z_t\}_{t \geq 0}$ is irreducible and non-reversible with a unique stationary distribu-
205 tion π' in which

$$\pi'(e_{ij}) = \pi_i P(i, j) = \pi_j P(j, i), \quad e_{ij} \in \mathcal{E}. \quad (46)$$

206 Also, for any scalar-valued function g , the asymptotic variance $\sigma_Z^2(g) \leq \sigma_X^2(g)$.

207 Now, we show how the non-Markov random walk with properties in Theorem F.1 can be included
208 in Lemma 3.1. For the original function $G : \mathbb{R}^d \times \mathcal{V} \rightarrow \mathbb{R}^d$, we define another function $\Phi :$
209 $\mathbb{R}^d \times \mathcal{E} \rightarrow \mathbb{R}^d$ such that $\Phi(\theta, e_{ij}) = G(\theta, j)$. Then, the SGD update rule (7) becomes $\theta_{t+1} =$
210 $\text{Proj}_{\Theta}(\theta_t - \gamma_{t+1} \nabla \Phi(\theta_t, Z_{t+1}))$. From (46), we have for any $\theta \in \mathbb{R}^d$,

$$\phi(\theta) \triangleq \mathbb{E}_{Z \sim \pi'} \Phi(\theta, Z) = \sum_{e_{ij} \in \mathcal{E}} \Phi(\theta, e_{ij}) \pi'(e_{ij}) = \sum_{i, j \in \mathcal{V}} G(\theta, j) \pi_j P(j, i) = \sum_{j \in \mathcal{V}} \frac{1}{n \pi_j} F(\theta, j) \pi_j = f(\theta), \quad (47)$$

211 showing the mean-field function $\phi(\theta)$ for $\Phi(\theta, Z)$ is the same as the objective function $f(\theta)$.

212 Next, we show assumptions (A1)-(A5) of Lemma 3.1 still hold for $\{Z_t\}_{t \geq 0}$ on the augmented state
 213 space \mathcal{E} and function Φ . Assumption (A1), (A5) and (A3) hold for function $\Phi(\theta, Z)$ because of our
 214 definition $\Phi(\theta, e_{ij}) = G(\theta, j)$. Assumption (A4) is satisfied because from Theorem F.1, $\{Z_t\}_{t \geq 0}$ is
 215 an irreducible and non-reversible Markov chain on the augmented state space \mathcal{E} and there always
 216 exists a solution (27) to (22). Assumption (A2) holds because matrix $\mathbf{K} = \nabla^2 \phi(\theta^*) = \nabla^2 f(\theta^*)$
 217 by (47). Therefore, we can say the Markov chain $\{Z_t\}_{t \geq 0}$ on the augmented state space \mathcal{E} , along
 218 with the newly defined function Φ , can apply Lemma 3.1. The asymptotic covariance matrix
 219 $\Sigma_Z \triangleq \Sigma_Z(\nabla \Phi(\theta^*, \cdot))$ is given as

$$\begin{aligned}
 \Sigma_Z &= \text{Var}_{\pi'}(\nabla \Phi(\theta^*, Z_0)) + \sum_{k \geq 1} \text{Cov}_{\pi'}(\nabla \Phi(\theta^*, Z_0), \nabla \Phi(\theta^*, Z_k)) \text{Cov}_{\pi'}(\nabla \Phi(\theta^*, Z_0), \nabla \Phi(\theta^*, Z_k))^T \\
 &= \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left\{ \left[\sum_{s=1}^t (\nabla \Phi(\theta^*, Z_s) - \mathbb{E}_{\pi'}(\nabla \Phi(\theta^*, \cdot))) \right] \left[\sum_{s=1}^t (\nabla \Phi(\theta^*, Z_s) - \mathbb{E}_{\pi'}(\nabla \Phi(\theta^*, \cdot))) \right]^T \right\} \\
 &= \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left\{ \left[\sum_{s=1}^t (\nabla G(\theta^*, Y_s) - \mathbb{E}_{\pi}(\nabla G(\theta^*, \cdot))) \right] \left[\sum_{s=1}^t (\nabla G(\theta^*, Y_s) - \mathbb{E}_{\pi}(\nabla G(\theta^*, \cdot))) \right]^T \right\} \\
 &= \Sigma_Y(\nabla G(\theta^*, \cdot)),
 \end{aligned} \tag{48}$$

220 where the third equality comes from (46) because

$$\mathbb{E}_{\pi'}[(\nabla \Phi(\theta^*, Z))] = \nabla f(\theta^*) = \sum_{j \in \mathcal{V}} \frac{1}{n\pi_j} \nabla F(\theta, j) \pi_j = \sum_{j \in \mathcal{V}} \pi_j \nabla G(\theta, j) = \mathbb{E}_{\pi}(\nabla G(\theta^*, \cdot)).$$

221 $\{Y_t\}_{t \geq 0}$ on the node space \mathcal{V} is the trajectory generated by $\{Z_t\}_{t \geq 0}$ on the augmented state space
 222 \mathcal{E} . Let \mathbf{V}_Z be the covariance matrix generated by the SGD algorithm driven by $\{Z_t\}_{t \geq 0}$. Denote
 223 $\Sigma_X \triangleq \Sigma_X(\nabla G(\theta^*, \cdot))$ the asymptotic covariance matrix and \mathbf{V}_X the covariance matrix in (20) from
 224 the original Markov chain $\{X_t\}$. Then, from Theorem F.1 we know the asymptotic variances of
 225 NBRW and SRW are ordered for any scalar-valued function. Then, with Theorem 3.6 (i) we know
 226 that the asymptotic covariance matrices of NBRW and SRW are Loewner ordered for any vector-
 227 valued function such that $\Sigma_Y(\nabla G(\theta^*, \cdot)) \leq_L \Sigma_X(\nabla G(\theta^*, \cdot))$. From (48) we have $\Sigma_Z \leq_L \Sigma_X$. By
 228 applying Theorem 3.6 (ii), we have $\mathbf{V}_Z \leq_L \mathbf{V}_X$.

229 G Proof of Lemma 4.2

230 Assume we have a vector-valued function $\mathbf{g} : [n] \rightarrow \mathbb{R}^d$. For shuffling without replacement, which
 231 traverses every node in each epoch with length n , we group all the terms in each k -th epoch (shown
 232 in Figure 3) and analyze the term $\sum_{i=(k-1)n}^{kn-1} (\mathbf{g}(X_i) - \mathbb{E}_{\pi}(\mathbf{g})) < \infty$ for $k \in \mathbb{Z}_+$. Note that we have

$$\sum_{i=(k-1)n}^{kn-1} \mathbf{g}(X_i) = \sum_{j=1}^n \mathbf{g}(j) \tag{49}$$

233 by definition of shuffling without replacement. By (49) and $\mathbb{E}_{\pi}(\mathbf{g}) = \frac{1}{n} \sum_{i=1}^n \mathbf{g}(i)$, we have

$$\sum_{i=(k-1)n}^{kn-1} (\mathbf{g}(X_i) - \mathbb{E}_{\pi}(\mathbf{g})) = \sum_{j=1}^n \mathbf{g}(j) - n \sum_{i=1}^n \frac{1}{n} \mathbf{g}(i) = 0. \tag{50}$$

234

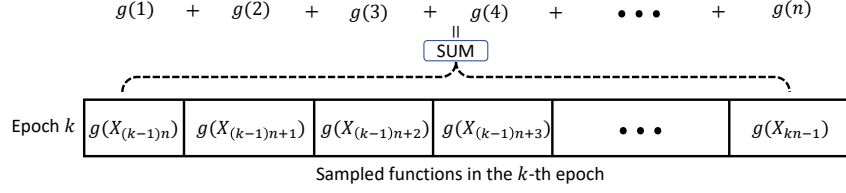


Figure 3: Diagram of the sampled functions in each epoch.

235 With the definition of the asymptotic covariance matrix, we have

$$\begin{aligned}
\Sigma_X(\mathbf{g}) &= \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left\{ \left(\sum_{i=1}^t (\mathbf{g}(X_s) - \mathbb{E}_\pi(\mathbf{g})) \right) \left(\sum_{i=1}^t (\mathbf{g}(X_s) - \mathbb{E}_\pi(\mathbf{g})) \right)^T \right\} \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left\{ \left[\sum_{i=1}^s (\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})) + \sum_{i=s+1}^t (\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})) \right] \right. \\
&\quad \cdot \left. \left[\sum_{i=1}^s (\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})) + \sum_{i=s+1}^t (\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})) \right]^T \right\} \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left\{ \left[\sum_{i=s+1}^t (\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})) \right] \left[\sum_{i=s+1}^t (\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})) \right]^T \right\},
\end{aligned} \tag{51}$$

where π is the uniform stationary distribution, $s \triangleq t - (t \bmod n)$ is the time at which the previous epoch ended before time t and we let $\sum_{i=t+1}^t (\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})) = 0$ by default. The third equality in (51) comes from (50). Note that there always exists a constant D such that $\|\mathbf{g}(i) - \mathbb{E}_\pi \mathbf{g}\|_2 < D$ for any $i \in [n]$ because of the boundedness of function \mathbf{g} . Then, we have for any t ,

$$\left\| \sum_{i=s+1}^t (\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})) \right\|_2 \leq \sum_{i=s+1}^t \|\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})\|_2 < (t - s - 1)D < nD < \infty,$$

where the second last inequality holds since $t - s - 1 = (t \bmod n) - 1 < n$. Back to (51), we have

$$\|\Sigma_X(\mathbf{g})\|_2 \leq \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[\left\| \sum_{i=s+1}^t (\mathbf{g}(X_i) - \mathbb{E}_\pi(\mathbf{g})) \right\|_2^2 \right] \leq \lim_{t \rightarrow \infty} \frac{nD}{t} = 0,$$

236 where the first inequality comes from Jensen's inequality. Finally, we have $\Sigma_X(\mathbf{g}) = \mathbf{0}$ such that the
237 asymptotic covariance matrices for both random and single shuffling are zero for any vector-valued
238 function \mathbf{g} .

239 H Proof of Proposition 4.3 and 4.4

240 H.1 Single Shuffling in SGD CLT Analysis

241 Single shuffling is seen as a time-homogeneous, irreducible, periodic Markov chain and we know
242 (A4) is the only requirement for Markov chain in the CLT result. As mentioned in [15] and [27]
243 Chapter 17, the necessary condition to ensure the existence of function \tilde{F} as in (27) is that the inverse
244 $(\mathbf{I} - \mathbf{P} + \mathbf{1}\pi^T)^{-1}$ exists. This is true for periodic Markov chain, and is shown in the following
245 lemma.

246 **Lemma H.1.** *The solution (27) to the Poisson equation (22) exists for an underlying finite, irreducible*
247 *periodic Markov chain with transition matrix $\mathbf{P} \in \mathbb{R}^{m \times m}$, stationary distribution π and period*
248 *$n \leq m$.*

249 *Proof.* From Perron–Frobenius theorem for irreducible, non-negative stochastic matrices [33], we
250 know there are n complex eigenvalues uniformly distributed on the unit circle, including the *unique*

251 eigenvalue with value 1. Other $m - n$ eigenvalues fall inside the unit circle, but still uniformly
 252 distributed on some circles with absolute value strictly smaller than 1 because transition matrix \mathbf{P} is
 253 similar to $e^{i\omega}\mathbf{P}$ where $i = \sqrt{-1}$ and $\omega = 2\pi/n$.

Denote $\lambda_1, \lambda_2, \dots, \lambda_m = 1$ the eigenvalues of the transition matrix \mathbf{P} and let \mathbf{J} be the Jordan norm
 form. There exists an invertible matrix $\mathbf{Q} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m]^T$ and $\mathbf{Q}^{-1} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m]$ such
 that $\mathbf{P} = \mathbf{Q}\mathbf{J}\mathbf{Q}^{-1}$ and $\mathbf{u}_i^T \mathbf{v}_i = 1$ for all $i \in [m]$. In particular, $\mathbf{u}_m = \boldsymbol{\pi}$ and $\mathbf{v}_m = \mathbf{1}$. Now,
 $(\lambda_m = 1, \mathbf{u}_m, \mathbf{v}_m)$ is also the PF eigenpair for the matrix $\boldsymbol{\Pi}$, which enables us to write down
 $\mathbf{1}\boldsymbol{\pi}^T = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^{-1}$, where $\boldsymbol{\Lambda} = \text{diag}(0, 0, \dots, 1)$. Therefore, $\mathbf{I} - \mathbf{P} + \mathbf{1}\boldsymbol{\pi}^T = \mathbf{Q}(\mathbf{I} - \mathbf{J} + \boldsymbol{\Lambda})\mathbf{Q}^{-1}$.
 Note that $\mathbf{I} - \mathbf{J} + \boldsymbol{\Lambda}$ is a new Jordan norm form with **non-zero** entries on the main diagonal. Assume
 $\mathbf{J}_i = \lambda_i \mathbf{I} + \mathbf{N}$ is one of its Jordan block with nilpotent matrix \mathbf{N} such that $\mathbf{N}^{p_i} = 0$ for some $p_i \geq 2$,
 then

$$\mathbf{J}_i^{-1} = \lambda_i^{-1}(\mathbf{I} + \lambda_i^{-1}\mathbf{N})^{-1} = \lambda_i^{-1}(\mathbf{I} - \lambda_i^{-1}\mathbf{N} + \dots + (\lambda_i)^{-p_i+1}\mathbf{N}^{p_i-1}),$$

254 showing that \mathbf{J}_i^{-1} exists for all Jordan blocks in the new Jordan norm form $\mathbf{I} - \mathbf{J} + \boldsymbol{\Lambda}$ because all λ_i
 255 are non-zero. Therefore, $(\mathbf{I} - \mathbf{P} + \mathbf{1}\boldsymbol{\pi}^T)^{-1}$ exists and (27) holds. \square

256 With Lemma H.1, single shuffling can indeed be included in the SGD CLT result, and its covariance
 257 matrix is the zero matrix $\mathbf{0}$ (from Lemma 4.2). Random shuffling is a *time-inhomogeneous* Markov
 258 chain due to its nature of reshuffling at the beginning of each epoch. Before providing our main
 259 proof, we first present the augmentation of the random shuffling sequence which transforms it into a
 260 *time-homogeneous* periodic Markov chain on the augmented state space.

261 H.2 Augmentation of Random Shuffling for CLT Analysis

262 By the definition of random shuffling, in each epoch of length n , the sampler traverses one permutation
 263 sequence drawn uniformly at random from the permutation sequence set with size $n!$. Due to its
 264 random nature across each epoch, random shuffling is not a Markov chain on state space $[n]$. In
 265 order to include random shuffling in the SGD CLT result, Lemma 3.1, we need to transform it into a
 266 Markov chain on an augmented state space.

267 Let $\{X_t\}_{t \geq 0}$ be the sequence generated by random shuffling. We first define an augmented
 268 state space \mathcal{S} , where for each state $s_t \triangleq \{\{A_j^{(t)}\}_{j \in [n]}, c_t\} \in \mathcal{S}$, the sequence $\{A_j^{(t)}\}_{j \in [n]} \triangleq$
 269 $\{X_{t-n+1}, X_{t-n+2}, \dots, X_t\}$ is of length n , and records the history of past n indices until time t . The
 270 integer $c_t \in \{1, 2, \dots, n\}$ is the time spent in current epoch at time t . For examples, consider the
 271 state space to be $\mathcal{S} = \{1, 2, \dots, 6\}$ in total and assume the sequence of visited states until $t = 8$ is
 272 $\{3, 6, 2, 1, 5, 4, 2, 5\}$. Here, $\{3, 6, 2, 1, 5, 4\}$ is one complete permutation sequence in the first epoch
 273 and $\{2, 5\}$ are in the second epoch. At time $t = 8$, the sampler is at index 5 and the sequence of past
 274 6 indices is $\{2, 1, 5, 4, 2, 5\}$ and $c_t = 2$, such that $s_8 = \{\{2, 1, 5, 4, 2, 5\}, 2\}$. In the next iteration
 275 $t = 9$, the sequence will be $\{1, 5, 4, 2, 5, X\}$ and $c_t = 3$, where X is the index that can be chosen
 276 from $\{1, 3, 4, 6\}$ uniformly at random because $\{2, 5\}$ have been chosen in the current epoch. Then,
 277 we have

$$s_9 = \begin{cases} \{\{1, 5, 4, 2, 5, 1\}, 3\} & \text{w.p } 1/4, \\ \{\{1, 5, 4, 2, 5, 3\}, 3\} & \text{w.p } 1/4, \\ \{\{1, 5, 4, 2, 5, 4\}, 3\} & \text{w.p } 1/4, \\ \{\{1, 5, 4, 2, 5, 6\}, 3\} & \text{w.p } 1/4. \end{cases} \quad (52)$$

278 Assume $s_{12} = \{2, 5, 6, 1, 3, 4, 6\}$ at $t = 12$, the next state $s_{13} = \{\{5, 6, 1, 3, 4, X\}, 1\}$ and X is
 279 chosen from $\{1, 2, \dots, 6\}$ uniformly at random.

280 Note that

- 281 • We only include *proper combination* of sequence $\{A_j\}_{j \in [n]}$ in the augmented state space,
 282 where ‘proper’ means the sequence is possible to appear with the current value of c_t .
 283 For instance, $\{\{2, 1, 5, 4, 2, 2\}, 2\}$ or $\{\{2, 1, 5, 1, 2, 5\}, 2\}$ is improper because $\{2, 2\}$ or
 284 $\{2, 1, 5, 1\}$ doesn’t exist in the permutation sequence in one epoch.
- 285 • Transition probability $P(s_t, s_{t+1})$ is possibly non-zero **only** when $c_{t+1} = c_t + 1$ for
 286 $c_t \leq n - 1$, or $c_{t+1} = 1$ when $c_t = n$.

287 Next, we show the proposition that will be used later to show that random shuffling can also be fitted
 288 into the CLT result.

289 **Proposition H.2.** $\{s_t\}_{t \geq 0}$ forms a finite, irreducible and periodic Markov chain with period n .

290 *Proof.* By our construction, the size of choice of $\{A_j\}_{j \in [n]}$ with $c = i$ is $(C_n^i)^2 i!(n-i)!$, because
 291 the first i indices has $C_n^i i!$ choices and remaining sequence has $C_n^{n-i} (n-i)!$ choices. The size of
 292 the augmented state space is $\sum_{i=1}^n (C_n^i)^2 i!(n-i)!$ and is still finite.

293 The *irreducibility* can be shown by $P(s_{t+2n} = s' | s_t = s) > 0$ because we can always construct two
 294 permutation sequences in two epochs; one including first i indices of $\{A'_j\}_{j \in [n]}$ in state s' and the
 295 other including the remaining sequences.

296 For *periodicity*, if the sequence $\{A_j^{(t)}\}_{j \in [n]}$ in the current state $s_t = s \in \mathcal{S}$ includes repeated in-
 297 dex at time t , e.g., index $i \in \{X_{t-n+1}, \dots, X_m\}$ (in the $\frac{m}{n}$ -th epoch) and $i \in \{X_{m+1}, \dots, X_t\}$
 298 (in the $(\frac{m}{n} + 1)$ -th epoch), then $i \notin \{X_{t+1}, \dots, X_{m+n}\}$ due to the nature of shuffling without
 299 replacement in an epoch, which leads to $P(s_{t+n} = s | s_t = s) = 0$. In addition, for $k \geq 2$ we can
 300 always construct intermediate sequences $\{X_{t+1}, \dots, X_{t+(k-1)n}\}$ such that $\{X_{t-n+1}, \dots, X_m\} =$
 301 $\{X_{t+(k-1)n+1}, \dots, X_{m+kn}\}$ and $\{X_{m+1}, \dots, X_t\} = \{X_{m+kn+1}, \dots, X_{t+kn}\}$, implying that
 302 $P(s_{t+kn} = s | s_t = s) > 0$ for $k \geq 2$. On the other hand, if $\{A_j^{(t)}\}_{j \in [n]}$ does not include repeated in-
 303 dex, $P(s_{t+kn} = s | s_t = s) > 0$ for $k \in \mathbb{N}$. We also note that $P(s_{t+j} = s | s_t = s) = 0$ for $s \in \mathcal{S}, j \neq$
 304 kn and $k \in \mathbb{N}$. Since $\{c_t\}$ by its definition is a periodic sequence $\{1, 2, \dots, n, 1, 2, \dots, n, 1, 2, \dots\}$
 305 with period of length n , we know that $c_t = c_{t+j}$ holds only when $j = kn$ for $k \in \mathbb{N}$. Then, for
 306 $j \neq kn$, we have $c_{t+j} \neq c_t$ such that $s_{t+j} \neq s_t$, which leads to $P(s_{t+j} = s | s_t = s) = 0$ for $j \neq kn$
 307 and $k \in \mathbb{N}$. Therefore, by definition of periodicity, the Markov chain is of period n . \square

308 Together with Lemma H.1 and Proposition H.2, we can see random shuffling can also be include in
 309 the SGD CLT.

310 H.3 Extension to Mini-batch Gradient Descent

311 Mini-batch gradient descent is another popular gradient descent variant and is widely used in the
 312 machine learning tools [10, 1, 31] to accelerate the learning process when compared to SGD. Instead
 313 of sampling a single element, mini-batch gradient descent samples multiple elements from $[n]$ in
 314 each iteration that form a batch.

315 To incorporate the notion of mini-batches in our SGD framework, we provide a reformulation of the
 316 general SGD iteration based on a similar formulation in [16] for the general analysis of SGD with
 317 *i.i.d* inputs. Consider a stochastic process $\{B_t\}_{t \geq 0}$ as the driving sequence, which randomly samples
 318 batches of size S (without replacement) from the state space $[n]$, that is $B_t \subset [n]$ and $|B_t| = S$ for all
 319 $t \geq 0$. Here we assume $[n] \bmod S = 0$ for simplicity. B_t will therefore refer to the batch chosen at
 320 any time $t > 0$. We assume that B_t for all $t > 0$ are *i.i.d* random variables drawn from a distribution
 321 \mathcal{P} , such that $\mathcal{P}(B) > 0$ is the probability with which a batch $B \subset [n]$ is picked. We associate with
 322 any batch B , $\mathbf{v}(B) \triangleq [\sum_{i \in B} \mathbf{e}_i] / \binom{N}{S} \mathcal{P}(B)$, where \mathbf{e}_i is the i 'th vector of the canonical basis of \mathbb{R}^d .
 323 We then denote $\mathbf{F}(\theta) \triangleq [F(\theta, 1), \dots, F(\theta, n)]^T$, and $\nabla \mathbf{F}(\theta) \triangleq [\nabla F(\theta, 1), \dots, \nabla F(\theta, n)]^T$ for all
 324 $\theta \in \Theta$. With this notation, we can rewrite the general update rule for mini-batch SGD as

$$\theta_{t+1} = \text{Proj}_{\Theta} (\theta_t - \gamma_{t+1} \nabla \mathbf{F}(\theta_t)^T \mathbf{v}(B_{t+1})). \quad (53)$$

325 Note that this way of defining the mini-batch based random input ensures that $\mathbb{E}_{\mathcal{P}}[\mathbf{F}(\theta)^T \mathbf{v}(\cdot)] = f(\theta)$
 326 for all $\theta \in \Theta$, maintaining the same objective function irrespective of the distribution from which
 327 batches are sampled.

328 With $X_t = B_t$ for all $t \geq 0$, and $\nabla G(\theta_t, X_{t+1}) = \nabla \mathbf{F}(\theta_t)^T \mathbf{v}(B_{t+1})$, the iteration (53) can still be
 329 written in the form of (7) with *i.i.d* input sequence $\{X_t\}_{t \geq 0}$. We can thus apply the CLT for SGD
 330 algorithms to the mini-batch SGD with *i.i.d* input, and in a similar fashion as (13) derive the explicit
 331 form of the asymptotic covariance matrix of (53), that is,

$$\Sigma_B(\nabla \mathbf{F}(\theta^*)^T \mathbf{v}(\cdot)) \triangleq \text{Var}_{B_0 \sim \mathcal{P}}(\nabla \mathbf{F}(\theta^*)^T \mathbf{v}(B_0)). \quad (54)$$

332 In practice, mini-batch gradient descent with shuffling is more widely used than *i.i.d* sampling [1], in
 333 which B_t is generated by shuffling-based method instead of independent drawn from a distribution.²
 334 At the beginning of each epoch, *Mini-batch gradient descent with random shuffling* shuffles the whole
 335 dataset $[n]$ and split it into small batches. On the other hand, *mini-batch gradient descent with single*
 336 *shuffling* only shuffles the dataset $[n]$ once before dividing it into batches, sticking to a predetermined
 337 sequence of batches for all epochs of the training process. As pointed out by [35], there is still a
 338 gap between practical implementation and theoretical analysis for mini-batch gradient descent with
 339 shuffling. Nevertheless, by extrapolating the analysis from Proposition 4.3, we are able to analyze the
 340 efficiency ordering of shuffling and *i.i.d* sampling in the mini-batch version, as stated next.

341 **Proposition H.3.** *Consider the mini-batch gradient descent (53) with stochastic inputs single/random*
 342 *shuffling $\{X_t\}_{t \geq 0}$ and *i.i.d* sampling $\{Y_t\}_{t \geq 0}$, we have $\theta_t^X, \theta_t^Y \xrightarrow[t \rightarrow \infty]{a.s.} \theta^*$ and $\mathbf{V}_X = \mathbf{0} \leq_L \mathbf{V}_Y$.*

343 *Proof.* Let $l \triangleq n/B \in \mathbb{N}$. We first give the following corollary.

344 **Corollary H.4.** *Mini-batch gradient descent with single shuffling is an irreducible, periodic Markov*
 345 *chain with period l .*

346 *Proof.* For single shuffling version, we divide the whole dataset $[n]$ into $B(1), B(2), \dots, B(l)$ and
 347 the corresponding sampling vector will be $\mathbf{v}(1), \mathbf{v}(2), \dots, \mathbf{v}(l)$. We shuffle the indices once, denoted
 348 by $a(1), a(2), \dots, a(l)$, and stick to this sequence all the time. Then, in each epoch, sampler will
 349 update θ_t according to the sequence $\mathbf{v}(a(1)), \mathbf{v}(a(2)), \dots, \mathbf{v}(a(l))$, where $\{\mathbf{v}_t\}_{t \geq 0}$ forms the finite,
 350 irreducible periodic Markov chain with period l . \square

351 Together with Corollary H.4 and Lemma H.1, mini-batch gradient descent with single shuffling can
 352 be included in the SGD CLT analysis and Theorem 3.6. Then, by Lemma H.1 and 4.2, mini-batch
 353 SGD with single shuffling can be applied to CLT result, which gives the asymptotic covariance matrix
 354 $\Sigma_{\mathbf{w}}(\nabla \mathbf{F}(\theta^*)^T \mathbf{v}(\cdot)) = \mathbf{0}$ and thus the covariance matrix in the CLT result is also zero.

355 For random shuffling version, we can use similar method as in Appendix H.2 to augment the state
 356 space, which forms the corollary as follows.

357 **Corollary H.5.** *$\{x_t\}_{t \geq 0}$ forms a finite, irreducible and periodic Markov chain with period l .*

358 *Proof.* Let \mathcal{X} be the augmented space, where state $x_t \triangleq \{\{W_j^{(t)}\}_{j \in [l]}, c_t\}$. Sequence $\{W_j^{(t)}\}_{j \in [l]} =$
 359 $\{\mathbf{v}_{t-l+1}, \mathbf{v}_{t-l+2}, \dots, \mathbf{v}_t\}$ records the last l selected batches and $c_t \in \{1, 2, \dots, l\}$ is the relative
 360 position of the batch in the current epoch at time t . The only difference for mini-batch version to the
 361 single element version is that we sample one batch of size B without replacement according to the
 362 indices yet to be chosen in the current epoch. Similar to the proof in Proposition H.2, $\{X_t\}_{t \geq 0}$ is
 363 also a finite, irreducible and periodic Markov chain with period l . \square

364 Corollary H.5 and Lemma H.1 show that mini-batch gradient descent with random shuffling can
 365 also be included in the SGD CLT analysis. However, we still need to check the form of asymptotic
 366 covariance matrix due to the augmentation. We follow the same idea from Appendix F and define
 367 a function $\Phi(\theta, x_t) \triangleq \mathbf{F}(\theta)^T \mathbf{v}(B_t)$. Then, from Lemma 4.2, asymptotic covariance matrix Σ_x is

²The reformulation (53) enables us to analyze mini-batch gradient descent with various stochastic processes that samples B_t , not just *i.i.d* input and shuffling. However, discussing general processes $\{B_t\}_{t \geq 0}$ is beyond the scope of this paper.

368 given as

$$\begin{aligned}
\Sigma_x &= \text{Var}_\pi(\nabla\Phi(\theta^*, x_0)) + \sum_{k \geq 1} \text{Cov}_\pi(\nabla\Phi(\theta^*, x_0), \nabla\Phi(\theta^*, x_k)) \text{Cov}_\pi(\nabla\Phi(\theta^*, x_0), \nabla\Phi(\theta^*, x_k))^T \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left\{ \left[\sum_{s=1}^t (\nabla\Phi(\theta^*, x_s) - \mathbb{E}_\pi(\nabla\Phi(\theta^*, \cdot))) \right] \left[\sum_{s=1}^t (\nabla\Phi(\theta^*, x_s) - \mathbb{E}_\pi(\nabla\Phi(\theta^*, \cdot))) \right]^T \right\} \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left\{ \left[\sum_{s=1}^t (\nabla\mathbf{F}(\theta^*)^T \mathbf{v}(B_s) - \nabla f(\theta^*)) \right] \left[\sum_{s=1}^t (\nabla\mathbf{F}(\theta^*)^T \mathbf{v}(B_s) - \nabla f(\theta^*)) \right]^T \right\} \\
&= \Sigma_x(\nabla\mathbf{F}(\theta^*)^T \mathbf{v}(\cdot)) = \mathbf{0},
\end{aligned} \tag{55}$$

369 where the third equality comes from the limiting distribution of random shuffling that is uniform.
370 Thus, the covariance matrix of random shuffling in the CLT result is also zero.

371 Above results show that both single shuffling and random shuffling in mini-batch SGD have higher
372 efficiency than mini-batch SGD with i.i.d sampling. \square

373 Proposition H.3 generalizes Proposition 4.3 (special case with mini-batch of size $S = 1$) in that the
374 same efficiency ordering between shuffling and *i.i.d* input holds true even with mini-batches.

375 I Simulation

376 In Appendix I.1, we give the details of our simulation setup for Figure 1, involving three reversible
377 Markov chains - the Metropolis-Hasting random walk (MHRW), a modification of MHRW (Modified-
378 MHRW) and fastest mixing Markov chain (FMMC), each having the uniform distribution as their
379 stationary measure. In Appendix I.2, we expand upon the numerical results in Section 5 by including
380 additional results for large graphs.

381 I.1 Details behind Figure 1

382 For the random walk SGD (RWSGD) simulation in Figure 1, we consider the problem of minimizing a
383 (scalar-valued) quadratic objective function

$$f(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n F(\theta, i) = \frac{1}{2n} \sum_{i=1}^n (\theta - b(i))^2, \tag{56}$$

384 where $\theta, b(i) \in \mathbb{R}$ for $i = 1, 2, \dots, n$ and n is the number of nodes on the graph. The minimizer is
385 given by $\theta^* \triangleq \arg \min_{\theta} f(\theta) = \frac{1}{n} \sum_{i=1}^n b(i)$. The RWSGD iteration for the objective function (56)
386 is then given by

$$\theta_{t+1} = \theta_t - \gamma_{t+1}(\theta_t - b(X_{t+1})), \tag{57}$$

387 where we choose $\gamma_t = 1/t^{0.9}$ and $\{X_t\}_{t \geq 0}$ is the stochastic input, e.g., MHRW, Modified-MHRW,
388 and FMMC.

389 In Figure 1, we simulate the SGD algorithm on two graphs; one is an 8-node graph \mathcal{G}_1 and the other
390 is a 5-node graph \mathcal{G}_2 . The two graphs are arbitrarily constructed while ensuring connectivity. See
391 Figure 4 for resulting topologies.

392 Now, we are ready to introduce the construction of three Markov chains on two graphs in Figure 4.

393 **MHRW:** Metropolis-Hasting algorithm [26] shows that the transition matrix of MHRW is constructed
394 in the following manner:

$$P(i, j) = \begin{cases} \min \left\{ \frac{1}{d_i}, \frac{1}{d_j} \right\}, & j \in N(i), \\ 1 - \sum_{j \in N(i)} P(i, j), & j = i, \end{cases} \tag{58}$$

395 where d_i is the degree of node i and $N(i)$ is the set of node i 's neighbors.

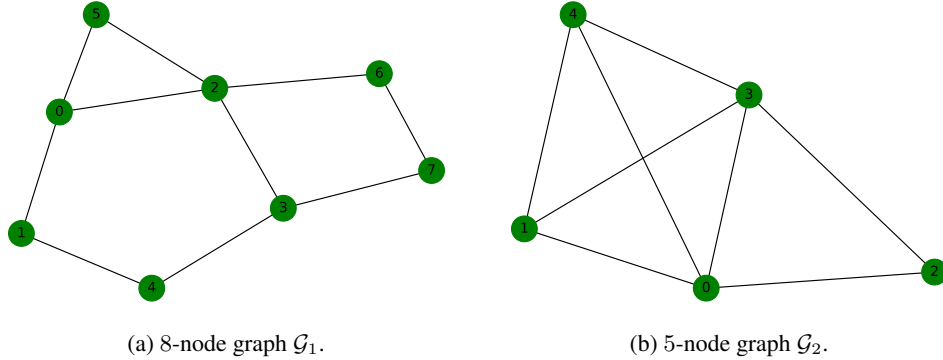


Figure 4: Topology of two graphs.

396 **Modified-MHRW:** To construct a ‘modified-MHRW’, which is more *efficient* than the standard
 397 MHRW,³ we employ the notion of ‘Peskun ordering’, originated from [32].

398 **Definition I.1** (Peskun ordering [32]). For two finite, ergodic, reversible Markov chains
 399 $\{X_t\}_{t \geq 0}, \{Y_t\}_{t \geq 0}$ on the state space \mathcal{V} with transition matrices $\mathbf{P}_X, \mathbf{P}_Y$ having the same station-
 400 ary distribution π , it is said that \mathbf{P}_Y dominates \mathbf{P}_X off the diagonal, written as $\mathbf{P}_X \preceq \mathbf{P}_Y$ if
 401 $P_X(i, j) \leq P_Y(i, j)$ for all $i, j \in \mathcal{V}$ and $i \neq j$.

402 We have the following lemma that connects the Peskun ordering to the efficiency ordering.

403 **Lemma I.2** ([32] Theorem 2.1.1). *If $\mathbf{P}_X \preceq \mathbf{P}_Y$, then $\sigma_X^2(g) \geq \sigma_Y^2(g)$ for any scalar-valued function*
 404 *g with $\mathbb{E}_\pi(g^2) < \infty$, that is, $\{Y_t\}_{t \geq 0}$ is more efficient than $\{X_t\}_{t \geq 0}$.*

405 We can manually construct a more efficient Markov chain by reducing the self-transition probability
 406 $P(i, i)$ of the MHRW and redistributing to off-diagonal entries, whenever possible, in a way that
 407 each row still sums to one and the resulting matrix is doubly-stochastic (i.e., the resulting Markov
 408 chain is reversible w.r.t the uniform distribution). In view of Lemma I.2, this modification improves
 409 the efficiency (smaller AV σ^2 compared to the standard MHRW).⁴

410 **FMMC:** FMMC is obtained by solving a semidefinite programming (proposed in problem (6) of [6]),
 411 which gives a Markov chain that minimizes the SLEM of the transition matrix over the entire class of
 412 reversible Markov chains w.r.t the uniform stationary distribution for a given graph topology. This is
 413 done numerically by using the CVXOPT package [12]. Later we will show in the simulation that
 414 FMMC indeed has the smallest SLEM compared to MHRW and Modified-MHRW.

415 In what follows, we index these three Markov chains with numbers in the subscript: MHRW (indexed
 416 by 1), Modified-MHRW (indexed by 2), and FMMC (indexed by 3). For graph \mathcal{G}_1 , the transition

³*Efficiency ordering* of Markov chains is introduced in Definition 3.5. In short, a Markov chain $\{X_t\}_{t \geq 0}$
 is more efficient than $\{Y_t\}_{t \geq 0}$ if the asymptotic variances (AV) satisfy $\sigma_X^2(g) \leq \sigma_Y^2(g)$ for any scalar-valued
 function g with $\mathbb{E}_\pi(g^2) < \infty$, where $\sigma_X^2(g)$ is defined in (4).

⁴Note that there can be many ways to modify the standard MHRW that make the Markov chain more efficient.
 The pursuit of the ‘optimal’ modification w.r.t the efficiency is out of the scope of this paper.

417 matrices of MHRW $\mathbf{P}_1^{\mathcal{G}_1}$, Modified-MHRW $\mathbf{P}_2^{\mathcal{G}_1}$, and FMMC $\mathbf{P}_3^{\mathcal{G}_1}$ are given by

$$\begin{aligned}
\mathbf{P}_1^{\mathcal{G}_1} &= \begin{bmatrix} 1/12 & 1/3 & 1/4 & 0 & 0 & 1/3 & 0 & 0 \\ 1/3 & 1/6 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 1/4 & 0 & 0 & 1/4 & 0 & 1/4 & 1/4 & 0 \\ 0 & 0 & 1/4 & 1/12 & 1/3 & 0 & 0 & 1/3 \\ 0 & 1/2 & 0 & 1/3 & 1/6 & 0 & 0 & 0 \\ 1/3 & 0 & 1/4 & 0 & 0 & 5/12 & 0 & 0 \\ 0 & 0 & 1/4 & 0 & 0 & 0 & 1/4 & 1/2 \\ 0 & 0 & 0 & 1/3 & 0 & 0 & 1/2 & 1/6 \end{bmatrix}, \\
\mathbf{P}_2^{\mathcal{G}_1} &= \begin{bmatrix} 0 & 0.35 & 0.25 & 0 & 0 & 0.4 & 0 & 0 \\ 0.35 & 0.02 & 0 & 0 & 0.63 & 0 & 0 & 0 \\ 0.25 & 0 & 0 & 0.25 & 0 & 0.25 & 0.25 & 0 \\ 0 & 0 & 0.25 & 0 & 0.37 & 0 & 0 & 0.38 \\ 0 & 0.63 & 0 & 0.37 & 0 & 0 & 0 & 0 \\ 0.4 & 0 & 0.25 & 0 & 0 & 0.35 & 0 & 0 \\ 0 & 0 & 0.25 & 0 & 0 & 0 & 0.13 & 0.62 \\ 0 & 0 & 0 & 0.38 & 0 & 0 & 0.62 & 0 \end{bmatrix}, \\
\mathbf{P}_3^{\mathcal{G}_1} &= \begin{bmatrix} 0.13 & 0.42 & 0.17 & 0 & 0 & 0.28 & 0 & 0 \\ 0.42 & 0.1 & 0 & 0 & 0.48 & 0 & 0 & 0 \\ 0.17 & 0 & 0 & 0.06 & 0 & 0.32 & 0.45 & 0 \\ 0 & 0 & 0.06 & 0.14 & 0.46 & 0 & 0 & 0.34 \\ 0 & 0.48 & 0 & 0.46 & 0.06 & 0 & 0 & 0 \\ 0.28 & 0 & 0.32 & 0 & 0 & 0.4 & 0 & 0 \\ 0 & 0 & 0.45 & 0 & 0 & 0 & 0.09 & 0.46 \\ 0 & 0 & 0 & 0.34 & 0 & 0 & 0.46 & 0.2 \end{bmatrix}.
\end{aligned} \tag{59}$$

418 For graph \mathcal{G}_2 , the transition matrices of MHRW $\mathbf{P}_1^{\mathcal{G}_2}$, Modified-MHRW $\mathbf{P}_2^{\mathcal{G}_2}$, and FMMC $\mathbf{P}_3^{\mathcal{G}_2}$ are
419 given by

$$\begin{aligned}
\mathbf{P}_1^{\mathcal{G}_2} &= \begin{bmatrix} 0 & 1/4 & 1/4 & 1/4 & 1/4 \\ 1/4 & 1/6 & 0 & 1/4 & 1/3 \\ 1/4 & 0 & 1/2 & 1/4 & 0 \\ 1/4 & 1/4 & 1/4 & 0 & 1/4 \\ 1/4 & 1/3 & 0 & 1/4 & 1/6 \end{bmatrix}, \\
\mathbf{P}_2^{\mathcal{G}_2} &= \begin{bmatrix} 0 & 0.25 & 0.25 & 0.25 & 0.25 \\ 0.25 & 0 & 0 & 0.25 & 0.5 \\ 0.25 & 0 & 0.5 & 0.25 & 0 \\ 0.25 & 0.25 & 0.25 & 0 & 0.25 \\ 0.25 & 0.5 & 0 & 0.25 & 0 \end{bmatrix}, \\
\mathbf{P}_3^{\mathcal{G}_2} &= \begin{bmatrix} 0.09 & 0.25 & 0.33 & 0.08 & 0.25 \\ 0.25 & 0.25 & 0 & 0.25 & 0.25 \\ 0.33 & 0 & 0.34 & 0.33 & 0 \\ 0.08 & 0.25 & 0.33 & 0.09 & 0.25 \\ 0.25 & 0.25 & 0 & 0.25 & 0.25 \end{bmatrix}.
\end{aligned} \tag{60}$$

420 In both (59) and (60), observe that Modified-MHRW and MHRW follow the Peskun ordering, i.e.,
421 $\mathbf{P}_1^{\mathcal{G}_1} \preceq \mathbf{P}_2^{\mathcal{G}_1}$ and $\mathbf{P}_1^{\mathcal{G}_2} \preceq \mathbf{P}_2^{\mathcal{G}_2}$, such that Modified-MHRW is more efficient than MHRW according to
422 Lemma I.2. In addition, the SLEMs of these matrices are given in Table 1, where FMMC has the
423 smallest SLEM in both graphs compared to MHRW and Modified-MHRW. Interestingly, Modified-
424 MHRW has larger SLEM than MHRW in graph \mathcal{G}_1 , which means Modified-MHRW can mix slower
425 than MHRW to the stationary distribution.

426 In Figure 5, we show the simulation result of each Markov chain in the RWSGD algorithm with
427 iteration (57) w.r.t MSE $\mathbb{E}\|\theta_t - \theta^*\|_2^2$ in graph \mathcal{G}_1 and \mathcal{G}_2 .⁵ In both graphs, Modified-MHRW (green
428 curve) performs better than MHRW (red curve) and FMMC (blue curve) with smallest MSE while it
429 has the largest SLEM shown in Table 1. This implies that the order of SLEM does not reflect the
430 order of MSE in the RWSGD algorithm.

⁵The reason we plot the same curves in each graph will be explained in the next paragraph.

	\mathcal{G}_1	\mathcal{G}_2
MHRW (β_1)	0.761	0.500
Modified-MHRW (β_2)	0.868	0.500
FMMC (β_3)	0.712	0.408

Table 1: SLEMs of the transition matrices in (59) and (60).

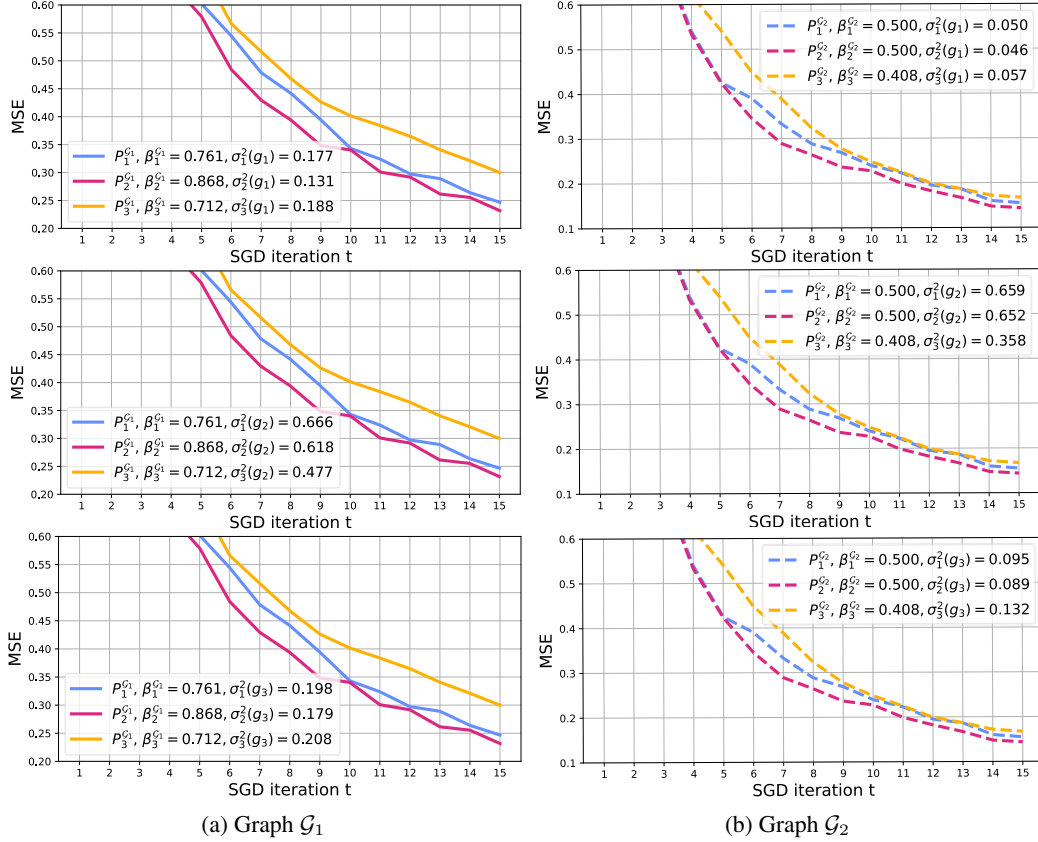


Figure 5: MSE $\mathbb{E} \|\theta_t - \theta^*\|_2^2$ of three Markov chains in the SGD algorithm with iteration (57).

431 In Figure 5, we repeat the plot in each graph three times with three different values of AVs inside the
432 legend, the reason being that we want to see if the performance of each Markov chain is related to the
433 AV $\sigma^2(g)$ and its test function g , other than SLEM solely. In the top row of Figure 5, as well as in
434 Figure 1, we choose the test function $g_1(i) = \nabla F(\theta^*, i)$ for $i = 1, 2, \dots, n$, where $\nabla F(\theta, i)$ is the
435 gradient of the local function $F(\theta, i)$ (56) w.r.t θ . In the middle row of Figure 5, the test function is
436 $g_2(i) = d_i$, which estimates the average degree of the graph. In the bottom row of Figure 5, the test
437 function is $g_3(i) = \mathbb{1}_{\{i=1\}}$, which estimates the probability of visiting node 1. We include the AVs of
438 all three test functions g_1, g_2, g_3 in the legend of Figure 5, e.g., $\sigma_3^2(g_1)$ is the AV of the test function
439 g_1 for FMMC.⁶ We observe that $\sigma_1^2(g_1) < \sigma_3^2(g_1)$ and $\sigma_1^2(g_3) < \sigma_3^2(g_3)$ while $\sigma_1^2(g_2) > \sigma_3^2(g_2)$
440 in both graphs. This means MHRW and FMMC are *not* efficiency ordered, which is possible because
441 efficiency ordering is a partial order such that not every two Markov chains can be ordered. On the
442 other hand, in both graphs, $\sigma_1^2(g_k) > \sigma_2^2(g_k)$ for $k = 1, 2, 3$, which is consistent with the fact that
443 the constructed Modified-MHRW is more efficient than MHRW. Regarding the MSE, we find that
444 Modified-MHRW performs better than MHRW in both graphs, which is in line with the efficiency
445 ordering. This leads us to conjecture that two efficiency ordered Markov chains might also have their
446 performance in the RWSGD algorithm ordered in the same way.

⁶The AV of the test function for each Markov chain is calculated by running a stochastic simulation for a long time and directly computing according to the definition in (4).

447 **Remark I.3.** Section 1.1 in [34] numerically compares the performance of a reversible Markov
448 chain and its non-reversible counterpart in the RWSGD algorithm w.r.t SLEM, and shows that the
449 non-reversible counterpart with smaller SLEM performs better. The main theorem therein is also
450 applicable to the comparison of two reversible Markov chains. However, as shown in Figure 5, we
451 provide examples to show that a reversible Markov chain with smaller SLEM does not necessarily
452 lead to smaller MSE in the RWSGD algorithm. Note that our results do not contradict the simulation
453 results in Section 1.1 of [34], since it is possible for a Markov chain to have both smaller AV and
454 smaller SLEM; which could be the case for the non-reversible counterpart in [34], although they
455 didn't specify the AV in their simulation. Moreover, their main theorem is an upper bound to the
456 error terms considered, which means that an SLEM-based ordering does not guarantee a performance
457 ordering of the error terms themselves, as also exemplified in Figure 5. All of these together imply
458 that SLEM alone cannot be the sole indicator of performance of the Markov chains as input sequences
459 for RWSGD algorithms. \square

460 I.2 Numerical Results on Large Graphs

461 We first specify the process of dataset generation for the sum-of-nonconvex functions $\hat{f}(\theta)$ in (14).
462 We generate random vectors $\mathbf{a}_1, \dots, \mathbf{a}_n, \mathbf{b} \in \mathbb{R}^{10}$ uniformly from $[0, 1]$ and ensure the invertibility
463 of $\sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^T$. Then, we randomly select half of the matrices in $\{\mathbf{D}_i\}_{i \in [n]}$ and assign $+1.1$ to their
464 j -th diagonal; other matrices are assigned -1.1 to j -th diagonal. We repeat the above process for all
465 diagonal values $j = 1, 2, \dots, 10$. This process guarantees $\sum_{i=1}^n \mathbf{D}_i = \mathbf{0}$.

466 We perform additional simulations on graph ‘AS-733’ [21] with 6474 nodes, and graph ‘wikiVote’
467 [20] with 889 nodes with the same objective functions $\hat{f}(\theta)$ and $\hat{f}(\theta)$ in (14). The simulation results
468 are given in Figure 6 and 7. We plot the curves of NBRW and SRW in the insets of Figures 6a,
469 6b, 7a, and 7b, with the same x,y axes but at linear scale, to better observe the difference in their
470 performance. For both objective functions, NBRW has smaller MSE than SRW and both random and
471 single shuffling perform better than uniform sampling, e.g., Figure 6a and 7a.⁷ This demonstrate that
472 NBRW and SRW are efficiency-ordered, which also holds for random/single shuffling and uniform
473 sampling. Note that since we simulate on large graphs, for the logistic regression problem, the
474 SGD algorithm with NBRW and SRW is yet to enter the asymptotic regime even in the 100,000-th
475 iteration, which can be explained by the blue and green increasing curves in the inset of Figure 6b and
476 7b. On the other hand, the curve of uniform sampling becomes flat and the curves of single/random
477 shuffling are starting to go down in Figure 6b, 6d and 7b, 7d, implying that they have entered the
478 asymptotic regime. These results are consistent to the observations in Figure 2, which support our
479 theory.

480 I.3 Additional Simulations on Non-convex Objective Function and SGD Variants

We first Regarding the SGD variants other than the vanilla SGD, central limit theorem (CLT) is less
well studied in the literature. To list a few, [19] studied variance reduced SGD (SVRG) and obtained
the CLT for constant step size. [2] analyzed Adam and their follow-up [3] extended the CLT for a
general SGD algorithm, which includes Stochastic Heavy Ball (SHB), Nesterov accelerated SGD
(NaSGD) and Adam. [23] established the CLT for momentum SGD (mSGD) and NaSGD under more
general conditions on the step size. However, all of these recent works focus only on the Martingale
difference noise

$$\mathbb{E}[\delta_{t+1} | \mathcal{F}_t \triangleq \sigma(\theta_0, X_0, X_1, \dots, X_t)] = \mathbb{E}[\nabla f(\theta_t) - \nabla F(\theta_t, X_{t+1}) | \mathcal{F}_t] = 0,$$

which is equivalent to saying that the input $\{X_t\}_{t \geq 0}$ is independently sampled from some identical
distribution for each time t (*i.i.d* input sequence). Meanwhile, for Markovian inputs,

$$\mathbb{E}[\delta_{t+1} | \mathcal{F}_t] = \sum_{i \in [n]} \pi_i \nabla F(\theta_t, i) - \sum_{i \in [n]} P(X_t, i) \nabla F(\theta_t, i) \neq 0$$

481 because $\pi_i \neq P(X_t, i)$ in general (unless $\{X_t\}_{t \geq 0}$ is an *i.i.d* sequence). It remains an open problem to
482 obtain the CLT for these SGD variants with general Markovian inputs, which would be a prerequisite

⁷The curves of NBRW and SRW in Figure 6a and 7a appear flat because they are plotted in the same figure
with uniform sampling and single/random shuffling, which have much smaller MSE. We plot the comparison
between NBRW and SRW separately in the inset.

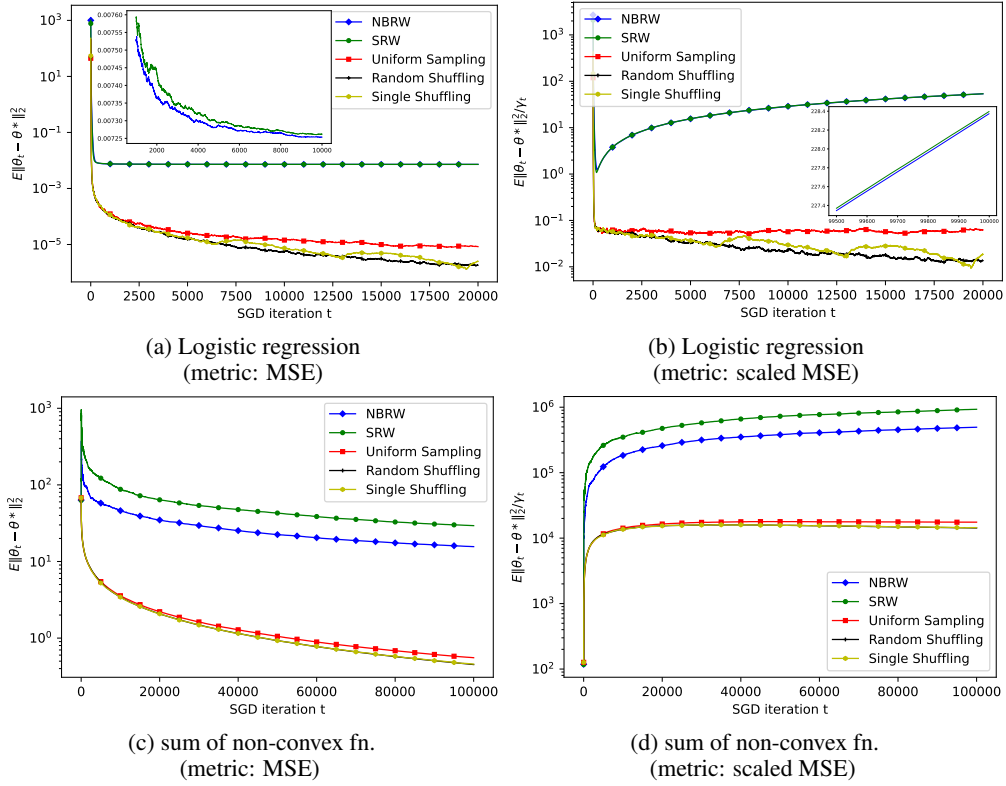


Figure 6: Performance comparison of different stochastic inputs on the graph ‘AS-733’.

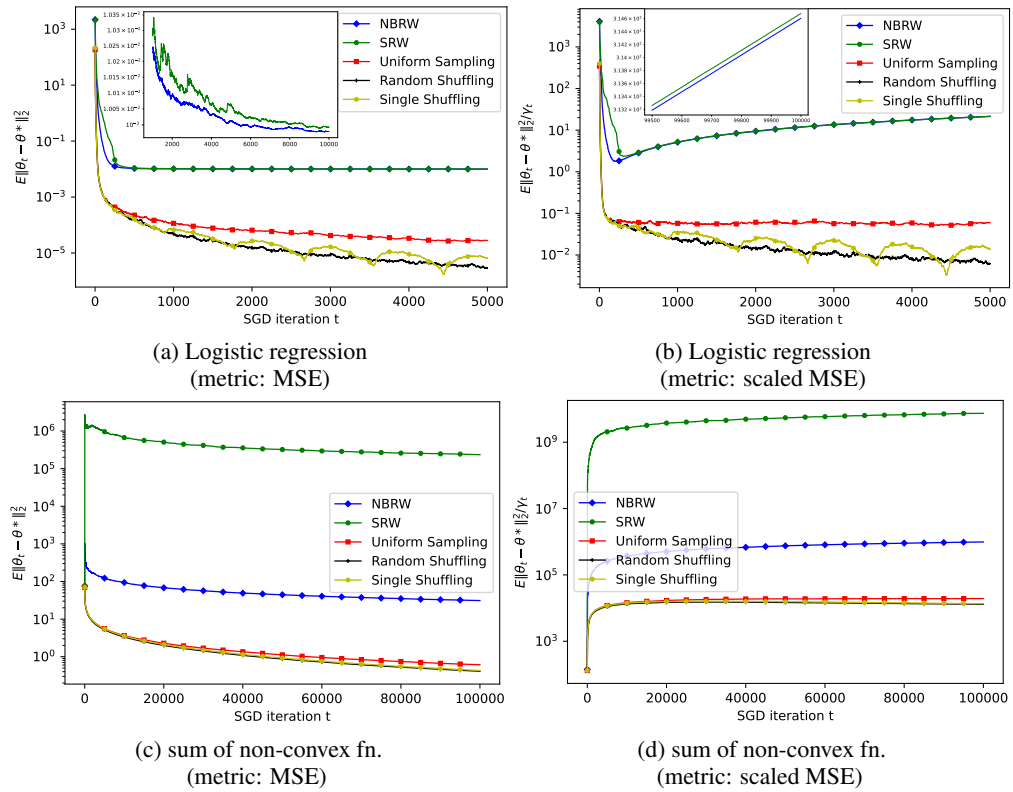


Figure 7: Performance comparison of different stochastic inputs on the graph ‘wikiVote’.

483 for our efficiency ordering. Indeed, one of our future works is to theoretically prove the CLT results for
 484 SGD variants with Markovian inputs and to carry over our efficiency ordering of different stochastic
 485 inputs.

486 Next, we will simulate two SGD variants, i.e., Nesterov accelerated SGD (NaSGD) and ADAM, on
 487 graph ‘‘AS-733’’ (as used in Appendix I.2) with two pair of stochastic inputs, i.e., NBRW versus SRW
 488 and shuffling methods versus *i.i.d* input sequence, with respect to both convex objective function and
 489 non-convex objective function. We choose the convex objective function $\hat{f}(\theta)$ from (14) such that

$$\hat{f}(\theta) = \frac{1}{n} \sum_{i=1}^n \theta^T (\mathbf{a}_i \mathbf{a}_i^T + \mathbf{D}_i) \theta + \mathbf{b}^T \theta, \quad (61)$$

490 where $\sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^T$ is invertible and $\sum_{i=1}^n \mathbf{D}_i = \mathbf{0}$. We can see $\nabla^2 \hat{f}(\theta) = \frac{2}{n} \sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^T$ is a positive
 491 semi-definite matrix and $\hat{f}(\theta)$ is convex. Then, we modify matrices $\{\mathbf{D}_i\}_{i \in [n]}$ such that the first
 492 element on the main diagonal of each matrix \mathbf{D}_i is subtracted by 0.1, and we denote the new matrices
 493 as $\{\mathbf{M}_i\}_{i \in [n]}$. We define a new function $\hat{g}(\theta)$ such that

$$\hat{g}(\theta) = \frac{1}{n} \sum_{i=1}^n \theta^T (\mathbf{a}_i \mathbf{a}_i^T + \mathbf{M}_i) \theta + \mathbf{b}^T \theta. \quad (62)$$

494 We numerically compute $\nabla^2 \hat{g}(\theta) = \frac{2}{n} \sum_{i=1}^n (\mathbf{a}_i \mathbf{a}_i^T + \mathbf{M}_i)$ and ensure it has at least one negative
 495 eigenvalue such that the objective function $\hat{g}(\theta)$ is non-convex. For Nesterov accelerated SGD, we
 496 employ the following iteration [23]

$$\begin{aligned} \theta_{t+1} &= u_t - \gamma_{t+1} \nabla G(u_t, X_{t+1}), \\ u_{t+1} &= \theta_{t+1} + \beta_{t+1} (\theta_{t+1} - \theta_t), \end{aligned} \quad (63)$$

497 where $\gamma_t = 1/0.9^t$ and $\beta_{t+1} \equiv \beta = 0.5$ in our settings. For ADAM, we use the following iteration
 498 [17]

$$\begin{aligned} g_{t+1} &= \nabla G(\theta_t, X_{t+1}), \\ m_{t+1} &= \alpha_1 m_t + (1 - \alpha_1) g_{t+1}, \\ v_{t+1} &= \alpha_2 v_t + (1 - \alpha_2) g_{t+1}^2, \\ m' &= m_{t+1} / (1 - \alpha_1^t), \\ v' &= v_{t+1} / (1 - \alpha_2^t), \\ \theta_{t+1} &= \theta_t - \gamma_t m' / (\sqrt{v'} + \epsilon), \end{aligned} \quad (64)$$

499 where $\gamma_t = 1/0.9^t$, $\alpha_1 = 0.9$, $\alpha_2 = 0.999$, $\epsilon = 10^{-8}$, g_{t+1}^2 is the element-wise square for the
 500 vector g_{t+1} and $\sqrt{v'}$ is the element-wise square root for the vector v' . In both (63) and (64),
 501 function $G(\theta, i) = \theta^T (\mathbf{a}_i \mathbf{a}_i^T + \mathbf{D}_i) \theta + \mathbf{b}^T \theta$ for convex objective function $\hat{f}(\theta)$ and $G(\theta, i) =$
 502 $\theta^T (\mathbf{a}_i \mathbf{a}_i^T + \mathbf{M}_i) \theta + \mathbf{b}^T \theta$ for non-convex objective function $\hat{g}(\theta)$.

503 The insets of Figure 8 are to enlarge the curves of NBRW and SRW in ADAM algorithm for the
 504 iteration $t \in [40000, 50000]$ to make them more distinguishable. In Figure 8, we show that for both
 505 convex and non-convex objective functions, the curves of NBRW are always below those of SRW
 506 in vanilla SGD, NaSGD and ADAM, respectively. This not only supports our theoretical results on
 507 vanilla SGD and both convex and non-convex objective functions, but also suggest that the efficiency
 508 ordering is still valid for other SGD variants.

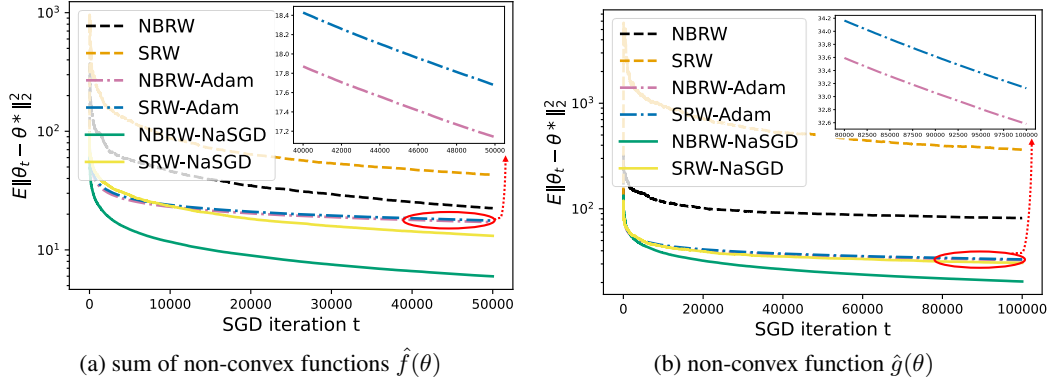


Figure 8: Performance comparison of NBRW and SRW in vanilla SGD, NaSGD and ADAM algorithms on the graph “AS-733”.

References

- 509
- 510 [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu
- 511 Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for
- 512 large-scale machine learning. In *12th USENIX symposium on operating systems design and*
- 513 *implementation (OSDI 16)*, pages 265–283, 2016.
- 514 [2] Anas Barakat and Pascal Bianchi. Convergence and dynamical behavior of the adam algorithm
- 515 for nonconvex stochastic optimization. *SIAM Journal on Optimization*, 31(1):244–274, 2021.
- 516 [3] Anas Barakat, Pascal Bianchi, Walid Hachem, and Sholom Schechtman. Stochastic optimization
- 517 with momentum: convergence, fluctuations, and traps avoidance. *Electronic Journal of Statistics*,
- 518 15(2):3892–3947, 2021.
- 519 [4] Albert Benveniste, Michel Métivier, and Pierre Priouret. *Adaptive algorithms and stochastic*
- 520 *approximations*, volume 22. Springer Science & Business Media, 2012.
- 521 [5] Vivek Borkar, Shuhang Chen, Adithya Devraj, Ioannis Kontoyiannis, and Sean Meyn. The ode
- 522 method for asymptotic statistics in stochastic approximation and reinforcement learning. *arXiv*
- 523 *preprint arXiv:2110.14427*, 2021.
- 524 [6] Stephen Boyd, Persi Diaconis, and Lin Xiao. Fastest mixing markov chain on a graph. *SIAM*
- 525 *review*, 46(4):667–689, 2004.
- 526 [7] Pierre Brémaud. *Markov chains: Gibbs fields, Monte Carlo simulation, and queues*, volume 31.
- 527 Springer Science & Business Media, 2013.
- 528 [8] VijaySekhar Chellaboina and Wassim M Haddad. *Nonlinear dynamical systems and control: A*
- 529 *Lyapunov-based approach*. Princeton University Press, 2008.
- 530 [9] Shuhang Chen, Adithya Devraj, Ana Busic, and Sean Meyn. Explicit mean-square error bounds
- 531 for monte-carlo and linear stochastic approximation. In *International Conference on Artificial*
- 532 *Intelligence and Statistics*, pages 4173–4183. PMLR, 2020.
- 533 [10] Francois Chollet et al. Keras, 2015.
- 534 [11] Bernard Delyon. Stochastic approximation with decreasing gain: Convergence and asymptotic
- 535 theory. Technical report, Université de Rennes, 2000.
- 536 [12] Steven Diamond and Stephen Boyd. Cvxpy: A python-embedded modeling language for convex
- 537 optimization. *The Journal of Machine Learning Research*, 17(1):2909–2913, 2016.
- 538 [13] Aymeric Dieuleveut, Alain Durmus, and Francis Bach. Bridging the gap between constant step
- 539 size stochastic gradient descent and markov chains. *The Annals of Statistics*, 48(3):1348–1382,
- 540 2020.

- 541 [14] Gersende Fort. Central limit theorems for stochastic approximation with controlled markov
542 chain dynamics. *ESAIM: Probability and Statistics*, 19:60–80, 2015.
- 543 [15] Peter W Glynn and Sean P Meyn. A liapounov bound for solutions of the poisson equation.
544 *The Annals of Probability*, pages 916–931, 1996.
- 545 [16] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter
546 Richtárik. Sgd: General analysis and improved rates. In *International Conference on Machine*
547 *Learning*, pages 5200–5209. PMLR, 2019.
- 548 [17] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*
549 *(Poster)*, 2015.
- 550 [18] Chul-Ho Lee, Xin Xu, and Do Young Eun. Beyond random walk and metropolis-hastings
551 samplers: why you should not backtrack for unbiased graph sampling. *ACM SIGMETRICS*
552 *Performance evaluation review*, 40(1):319–330, 2012.
- 553 [19] Jinlong Lei and Uday V Shanbhag. Variance-reduced accelerated first-order methods: Central
554 limit theorems and confidence statements. *arXiv preprint arXiv:2006.07769*, 2020.
- 555 [20] Jure Leskovec, Daniel Huttenlocher, and Jon Kleinberg. Signed networks in social media.
556 In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages
557 1361–1370, 2010.
- 558 [21] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graphs over time: densification laws,
559 shrinking diameters and possible explanations. In *Proceedings of the eleventh ACM SIGKDD*
560 *international conference on Knowledge discovery in data mining*, pages 177–187, 2005.
- 561 [22] Rong-Hua Li, Jeffrey Xu Yu, Lu Qin, Rui Mao, and Tan Jin. On random walk based graph
562 sampling. In *2015 IEEE 31st International Conference on Data Engineering (ICDE)*, pages
563 927–938. IEEE, 2015.
- 564 [23] Tiejun Li, Tiannan Xiao, and Guoguo Yang. Revisiting the central limit theorems for the
565 sgd-type methods. *arXiv preprint arXiv:2207.11755*, 2022.
- 566 [24] Yongkun Li, Zhiyong Wu, Shuai Lin, Hong Xie, Min Lv, Yinlong Xu, and John CS Lui. Walking
567 with perception: Efficient random walk sampling via common neighbor awareness. In *2019*
568 *IEEE 35th International Conference on Data Engineering (ICDE)*, pages 962–973. IEEE, 2019.
- 569 [25] Daniel Liberzon. *Calculus of variations and optimal control theory*. Princeton university press,
570 2011.
- 571 [26] Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and
572 Edward Teller. Equation of state calculations by fast computing machines. *The journal of*
573 *chemical physics*, 21(6):1087–1092, 1953.
- 574 [27] Sean P Meyn and Richard L Tweedie. *Markov chains and stochastic stability*. Springer Science
575 & Business Media, 2012.
- 576 [28] Sujit Kumar Mitra, P Bhimasankaram, and Saroj B Malik. *Matrix partial orders, shorted*
577 *operators and applications*, volume 10. World Scientific, 2010.
- 578 [29] Wenlong Mou, Chris Junchi Li, Martin J Wainwright, Peter L Bartlett, and Michael I Jordan. On
579 linear stochastic approximation: Fine-grained polyak-ruppert and non-asymptotic concentration.
580 In *Conference on Learning Theory*, pages 2947–2997. PMLR, 2020.
- 581 [30] Radford M Neal. Improving asymptotic variance of mcmc estimators: Non-reversible chains
582 are better. Technical report, Department of Statistics, University of Toronto, July 2004.
- 583 [31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
584 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
585 style, high-performance deep learning library. *Advances in neural information processing*
586 *systems*, 32:8026–8037, 2019.

- 587 [32] Peter H Peskun. Optimum monte-carlo sampling using markov chains. *Biometrika*, 60(3):607–
588 612, 1973.
- 589 [33] Eugene Seneta. *Non-negative matrices and Markov chains*. Springer Science & Business Media,
590 2006.
- 591 [34] Tao Sun, Yuejiao Sun, and Wotao Yin. On markov chain gradient descent. In *Proceedings of the*
592 *32nd International Conference on Neural Information Processing Systems*, pages 9918–9927,
593 2018.
- 594 [35] Chulhee Yun, Shashank Rajput, and Suvrit Sra. Minibatch vs local sgd with shuffling: Tight
595 convergence bounds and beyond. *arXiv preprint arXiv:2110.10342*, To appear in *ICLR 2022*,
596 2021.