

# [Fall 2025] ECE 5290/7290 and ORIE 5290 Distributed Optimization for Machine Learning and AI

#### Homework 3

Gradescope Due: October 20th at 11:59PM

## Objective of This Assignment

The goal of this assignment is to deepen your understanding of distributed optimization methods used in machine learning. You will analyze the convergence of consensus averaging (synchronous and gossip), understand how spectral properties control rates, compare mini-batch and parallel SGD, and study communication—computation trade-offs in local SGD. A coding problem will help you visualize convergence behaviors and the impact of network topology.

### Instruction of Homework Submission

This assignment includes both an analytical part and a coding part (Problem 5). We will use Gradescope to check the correctness of your code. Therefore, you will see two separate assignments on Gradescope.

- (a) A starter .py file is provided for Problem 5. Do not change the function names, signatures, or filename.
- (b) Upload the written PDF to **Homework 3** and the code to **Homework 3 Coding**.

## Question 1: Consensus Averaging and Spectral Gap (20 points)

Let  $\mathbf{x}(0) \in \mathbb{R}^N$  be scalar values held by N nodes on an undirected connected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ . The synchronous (Jacobi) consensus update is

$$\mathbf{x}(t+1) = \mathbf{W} \mathbf{x}(t), \qquad \mathbf{W} = \mathbf{I} - \alpha \mathbf{L},$$

where **L** is the graph Laplacian matrix and  $0 < \alpha < 1/\lambda_{\max}(\mathbf{L}); \bar{x} = \frac{1}{N}\mathbf{1}^{\top}\mathbf{x}(0)$  is the global average that the algorithm wants to achieve the consensus on. Define the consensus error as  $\mathbf{z}(t) = \mathbf{x}(t) - \bar{x}\mathbf{1}$ .

(a) (5 points) Show that 1 is an eigenvector of W with eigenvalue 1; further, W is symmetric and doubly stochastic for the chosen  $\alpha$ .

- (1 Point) L1 = 0 implies W1 =  $(I \alpha L)1 = 1$ .
- (2 Point) Since L is symmetric, so is W.
- (2 Point) Row sums: W1 = 1. Symmetry  $\Rightarrow$  column sums also 1 (doubly stochastic).

(b) (5 points) Let 
$$1 = \lambda_1(\mathbf{W}) \ge \lambda_2(\mathbf{W}) \ge \cdots \ge \lambda_N(\mathbf{W}) > -1$$
. Prove the following  $\|\mathbf{z}(t)\|_2 \le |\lambda_2(\mathbf{W})|^t \|\mathbf{z}(0)\|_2$ .

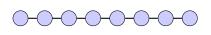
(3 Points) Decompose  $\mathbf{z}(0)$  in  $\mathbf{W}$  as orthonormal eigenbasis  $\{\mathbf{u}_i\}$  with  $\mathbf{u}_1 = \frac{1}{\sqrt{N}}\mathbf{1}$ . Notice that  $\langle \mathbf{z}(0), \mathbf{u}_1 \rangle = 0$ .

(2 Points) Therefore, we have  $\mathbf{z}(t) = \sum_{i=2}^{N} \lambda_i^t \langle \mathbf{z}(0), \mathbf{u}_i \rangle \mathbf{u}_i$ , so  $\|\mathbf{z}(t)\|_2 \leq \max_{i \geq 2} |\lambda_i|^t \|\mathbf{z}(0)\|_2$ .

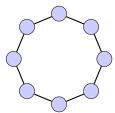
### (c) (10 points) Consider three graphs with N nodes and $W = I - \alpha L$ :

- Path graph  $P_N$  with  $\alpha = 1/\Delta$  (max degree  $\Delta = 2$ ).
- Ring graph  $C_N$  with  $\alpha = 1/\Delta$  (max degree  $\Delta = 2$ ).
- Complete graph  $K_N$  with  $\alpha = 1/N$ .

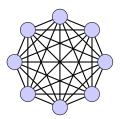
Recall  $\lambda_2(\mathbf{L})$  for each graph from class and give a clean upper bound on  $|\lambda_2(\mathbf{W})| = |1 - \alpha \lambda_2(\mathbf{L})|$ . Which graph mixes fastest?



Path graph  $P_N$ .



Ring Graph -  $C_N$ 



Complete Graph -  $\mathsf{C}_N$ 

(2 Points) Known spectra:  $\mathbf{L}(\mathsf{P}_N)$  has  $\lambda_2 \asymp c/N^2$ ;  $\mathbf{L}(\mathsf{C}_N)$  has  $\lambda_2 = 2(1 - \cos(2\pi/N)) \asymp c'/N^2$ ;  $\mathbf{L}(\mathsf{K}_N)$  has  $\lambda_2 = N$ .

(3 Points) Then  $|\lambda_2(\mathbf{W})| = |1 - \alpha \lambda_2(\mathbf{L})|$ . With  $\alpha = 1/2$  for path/ring,  $|\lambda_2(\mathbf{W})| \approx 1 - c/N^2$ ; with  $\alpha = 1/N$  on complete graph,  $|\lambda_2(\mathbf{W})| = |1 - (1/N) \cdot N| = 0$  for the second eigenvalue, hence one-step consensus in the idealized model. Thus  $\mathsf{K}_N$  mixes fastest, then  $\mathsf{C}_N \sim \mathsf{P}_N$  (both  $O(N^2)$  time to accuracy).

# Question 2: Randomized Gossip (Pairwise Averaging) (20 points)

At each step t, pick an edge  $(i, j) \in \mathcal{E}$  uniformly at random; the two endpoints average their values:

$$x_i(t+1) = x_j(t+1) = \frac{1}{2}(x_i(t) + x_j(t)), \quad x_\ell(t+1) = x_\ell(t) \text{ for } \ell \notin \{i, j\}.$$

Define the disagreement potential  $V(t) = \sum_{m=1}^{N} (x_m(t) - \bar{x})^2$  with  $\bar{x} = \frac{1}{N} \mathbf{1}^{\top} \mathbf{x}(0)$  as the global average.

(a) (8 points) Show that V(t) is nonincreasing and derive

$$\mathbb{E}\big[V(t+1) \mid \mathbf{x}(t)\big] = V(t) - \frac{1}{2|\mathcal{E}|} \sum_{(i,j)\in\mathcal{E}} (x_i(t) - x_j(t))^2.$$

2

Only the selected pair (i,j) changes. One checks  $V(t+1)-V(t)=-\frac{1}{2}\big(x_i-x_j\big)^2$ . Taking expectation over the uniform edge choice yields the identity.

(b) (6 points) Use  $\sum_{(i,j)\in\mathcal{E}} (x_i - x_j)^2 = \mathbf{x}^\top \mathbf{L} \mathbf{x} \ge \lambda_2(\mathbf{L}) \|\mathbf{x} - \bar{x}\mathbf{1}\|^2 = \lambda_2(\mathbf{L}) V$  to prove

$$\mathbb{E}\big[V(t+1) \mid \mathbf{x}(t)\big] \leq \left(1 - \frac{\lambda_2(\mathbf{L})}{2|\mathcal{E}|}\right) V(t).$$

Combine part (a) with the spectral lower bound to get the linear contraction in conditional expectation.

(c) (6 points) For the complete graph with uniform edge sampling, show

$$\mathbb{E}\big[V(t)\big] \leq \left(1 - \frac{1}{N}\right)^t V(0).$$

(Hint:  $\lambda_2(\mathbf{L}_{\mathsf{K}_N}) = N$  and  $|\mathcal{E}| = \frac{N(N-1)}{2}$ .)

Plug  $\lambda_2(\mathbf{L})=N$  and  $|\mathcal{E}|=\frac{N(N-1)}{2}$  into part (b):  $1-\frac{\lambda_2}{2|\mathcal{E}|}=1-\frac{N}{N(N-1)}=1-\frac{1}{N-1}\leq 1-\frac{1}{N}$ , yielding the stated bound (or the slightly sharper  $1-\frac{1}{N-1}$ ).

# Question 3: Mini-batch vs. Parallel SGD (15 points)

Consider empirical risk  $F(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} f_i(\mathbf{x})$ . Let K workers each draw independent mini-batches of size B (with replacement). Two **commonly used yet alternative** updates at iteration t:

1) (Option 1: Single-node mini-batch) One node samples a single mini-batch  $\mathcal{B}$  of size KB:

$$g(\mathbf{x}(t)) = \frac{1}{KB} \sum_{i \in \mathcal{B}} \nabla f_i(\mathbf{x}(t)), \text{ update } \mathbf{x}(t+1) = \mathbf{x}(t) - \eta g(\mathbf{x}(t)).$$

2) (Option 2: Parallel averaging) Each worker computes  $g^{(k)}(\mathbf{x}(t)) = \frac{1}{B} \sum_{i \in \mathcal{B}_k} \nabla f_i(\mathbf{x}(t))$ , then average:

$$\bar{g}(\mathbf{x}(t)) = \frac{1}{K} \sum_{k=1}^{K} g^{(k)}(\mathbf{x}(t)) \text{ update } \mathbf{x}(t+1) = \mathbf{x}(t) - \eta \, \bar{g}(\mathbf{x}(t)).$$

(a) (5 points) Show both estimators in Option 1 and Option 2 are unbiased for  $\nabla F(\mathbf{x}(t))$ .

By linearity of expectation and i.i.d. sampling with replacement,  $\mathbb{E}[g^{(k)}] = \nabla F$ , so  $\mathbb{E}[\bar{g}] = \nabla F$ . Similarly  $\mathbb{E}[g] = \nabla F$ .

(b) (5 points) Assuming the per-sample gradient variance is bounded by  $\mathbb{E}\|\nabla f_i(\mathbf{x}) - \nabla F(\mathbf{x})\|^2 \leq \sigma^2$ , prove

$$\operatorname{Var}(\bar{g}) = \frac{\sigma^2}{KB}, \quad \operatorname{Var}(g) = \frac{\sigma^2}{KB}.$$

Independence across workers and across samples yields variance additivity: each worker has variance  $\sigma^2/B$ , averaging K i.i.d. workers gives  $\sigma^2/(KB)$ . The single-node KB-batch has the same variance reduction.

(c) (5 points) Discuss when parallel averaging and single-node mini-batch are not equivalent in practice.

They differ under: (i) communication latency/stragglers (parallel cost per step), (ii) data heterogeneity across workers (breaks identical sampling), (iii) asynchronous implementations (stale gradients), and (iv) systems constraints (bandwidth, memory layout).

## Question 4: Local SGD - Communication/Computation Trade-offs (25 points)

Consider K workers collaboratively minimizing a global objective

$$F(\mathbf{x}) = \frac{1}{K} \sum_{k=1}^{K} F_k(\mathbf{x})$$
 with  $F_k(\mathbf{x}) = \frac{1}{B} \sum_{i \in \mathcal{B}_k} f_i(\mathbf{x})$ ,

where worker k has access to local data and computes stochastic gradients  $g^{(k)}$ . Each worker starts from a common model  $\mathbf{x}(t)$ , performs  $\tau$  local SGD steps (with  $\mathbf{x}_0^{(k)}(t) = \mathbf{x}(t)$ ):

$$\mathbf{x}_{s+1}^{(k)}(t) = \mathbf{x}_{s}^{(k)}(t) - \eta g_{s}^{(k)}(t), \qquad s = 0, 1, \dots, \tau - 1,$$

and then all workers synchronize by averaging:

$$\mathbf{x}(t+1) = \frac{1}{K} \sum_{k=1}^{K} \mathbf{x}_{\tau}^{(k)}(t).$$

Assume that each local function  $F_k$  is  $\mu$ -strongly convex and L-smooth, and that the stochastic gradients have bounded variance  $\sigma^2$ .

(a) (15 points) For ECE 7290 students: (Sketch) Show that compared to fully synchronized minibatch SGD (from Question 3), local SGD includes an additional *drift term* due to model divergence between averaging rounds. Derive a bound of the form ( $\eta^2$  in the initial version should be  $\eta$ )

$$\mathbb{E}\|\mathbf{x}(t+1) - \mathbf{x}^\star\|^2 \ \leq \ \rho^\tau \, \mathbb{E}\|\mathbf{x}(t) - \mathbf{x}^\star\|^2 \ + \ C_1 \, \frac{\eta \sigma^2}{\mu K} \ + \ C_2 \, \frac{\eta \tau}{\eta \tau} \Gamma^2,$$

where  $\Gamma^2$  captures the gradient dissimilarity (data heterogeneity) across nodes. Explain the dependence on  $\tau$ . For homogeneous data ( $\Gamma^2 \approx 0$ ), what  $\tau$  do you recommend? For heterogeneous data (large  $\Gamma^2$ ), how would you adjust  $\tau$ ?

The first term is geometric contraction under strong convexity. The second is the steady-state noise floor, reduced by K. The third arises from drift: local iterates deviate across nodes by  $O(\eta\tau)$ , and heterogeneity scales this to a bias  $O(\eta^2\tau\Gamma^2)$ . Larger  $\tau$  reduces communication but increases drift.

**Assumption 1 (Heterogeneity)** There exists  $\Gamma^2 \geq 0$  such that for all  $\mathbf{x}$ ,

$$\frac{1}{K} \sum_{k=1}^{K} \left\| \nabla F_k(\mathbf{x}) - \nabla F(\mathbf{x}) \right\|^2 \leq \Gamma^2, \qquad F(\mathbf{x}) = \frac{1}{K} \sum_{k=1}^{K} F_k(\mathbf{x}).$$

Define the per-step average iterate and the decomposition terms

$$\bar{\mathbf{x}}_s(t) := \frac{1}{K} \sum_{k=1}^K \mathbf{x}_s^{(k)}(t), \quad \zeta_s(t) := \frac{1}{K} \sum_{k=1}^K \left( g_s^{(k)}(t) - \nabla F_k(\mathbf{x}_s^{(k)}(t)) \right),$$

$$\Delta_s(t) := \frac{1}{K} \sum_{k=1}^K \left( \nabla F_k(\mathbf{x}_s^{(k)}(t)) - \nabla F_k(\bar{\mathbf{x}}_s(t)) \right).$$

Then the averaged update is

$$\bar{\mathbf{x}}_{s+1}(t) = \bar{\mathbf{x}}_s(t) - \eta \Big( \nabla F(\bar{\mathbf{x}}_s(t)) + \Delta_s(t) + \zeta_s(t) \Big).$$

Lemma 1 (One-step recursion of averaged iterate) If  $\eta \leq 1/L$ , then for all s,

$$\mathbb{E}\left\|\bar{\mathbf{x}}_{s+1}(t) - \mathbf{x}^{\star}\right\|^{2} \leq (1 - \eta\mu) \,\mathbb{E}\left\|\bar{\mathbf{x}}_{s}(t) - \mathbf{x}^{\star}\right\|^{2} + \eta^{2} \,\mathbb{E}\left\|\zeta_{s}(t)\right\|^{2} + \left(\eta^{2} + \frac{\eta}{\mu}\right) \,\mathbb{E}\left\|\Delta_{s}(t)\right\|^{2}. \tag{1}$$

Let  $oldsymbol{e}_s(t) := ar{\mathbf{x}}_s(t) - \mathbf{x}^\star.$  From the averaged update,

$$e_{s+1} = e_s - \eta (\nabla F(\bar{\mathbf{x}}_s) + \Delta_s + \zeta_s).$$

Expand, condition on  $\mathcal{F}_s$ , and use  $\mathbb{E}[\zeta_s \mid \mathcal{F}_s] = 0$  to remove mixed terms with  $\zeta_s$ :

$$\mathbb{E}\left[\|\boldsymbol{e}_{s+1}\|^2 \mid \mathcal{F}_s\right] = \left\|\boldsymbol{e}_s - \eta \nabla F(\bar{\mathbf{x}}_s)\right\|^2 + \eta^2 \|\Delta_s\|^2 + \eta^2 \mathbb{E}\left[\|\zeta_s\|^2 \mid \mathcal{F}_s\right] + 2\eta \langle \boldsymbol{e}_s - \eta \nabla F(\bar{\mathbf{x}}_s), -\Delta_s \rangle.$$

Upper bound the inner product by Young's inequality with parameter  $\mu > 0$ :

$$2\eta \langle \boldsymbol{e}_s - \eta \nabla F(\bar{\mathbf{x}}_s), -\Delta_s \rangle \leq \eta \mu \|\boldsymbol{e}_s - \eta \nabla F(\bar{\mathbf{x}}_s)\|^2 + \frac{\eta}{\mu} \|\Delta_s\|^2.$$

Hence

$$\mathbb{E}\big[\|\boldsymbol{e}_{s+1}\|^2\mid\mathcal{F}_s\big] \;\leq\; \big(1+\eta\mu\big)\,\big\|\boldsymbol{e}_s-\eta\nabla F(\bar{\mathbf{x}}_s)\big\|^2+\Big(\eta^2+\frac{\eta}{\mu}\Big)\,\|\Delta_s\|^2+\eta^2\mathbb{E}\big[\|\zeta_s\|^2\mid\mathcal{F}_s\big].$$

Using the standard GD contraction for  $\mu$ -strongly convex, L-smooth F with  $\eta \leq 1/L$ ,

$$\|\boldsymbol{e}_s - \eta \nabla F(\bar{\mathbf{x}}_s)\|^2 \le (1 - \eta \mu) \|\boldsymbol{e}_s\|^2,$$

we obtain

$$\mathbb{E} \big[ \| \boldsymbol{e}_{s+1} \|^2 \mid \mathcal{F}_s \big] \ \leq \ \underbrace{(1 + \eta \mu)(1 - \eta \mu)}_{\leq \ 1 - \eta \mu \ \text{for} \ \eta \mu \leq 1/2} \| \boldsymbol{e}_s \|^2 + \left( \eta^2 + \frac{\eta}{\mu} \right) \| \Delta_s \|^2 + \eta^2 \mathbb{E} \big[ \| \zeta_s \|^2 \mid \mathcal{F}_s \big].$$

Taking the total expectation yields (1).

**Lemma 2 (2nd moment of noise** + **drift)** With L-smooth,  $\mu$ -strongly convex and Assumption 1,

$$\mathbb{E}\|\Delta_s(t) + \zeta_s(t)\|^2 \le 2L^2 D_s(t) + 2\Gamma^2 + \frac{2\sigma^2}{K}, \qquad D_s(t) := \frac{1}{K} \sum_{k=1}^K \mathbb{E}\|\mathbf{x}_s^{(k)}(t) - \bar{\mathbf{x}}_s(t)\|^2.$$
 (2)

Sketch. By  $(a+b)^2 \le 2a^2 + 2b^2$  and the variance bound on  $\zeta_s$ ,  $\mathbb{E}\|\zeta_s\|^2 \le \sigma^2/K$ . Split  $\Delta_s = A_s + B_s$  with

$$A_s := \frac{1}{K} \sum_k \left( \nabla F_k(\mathbf{x}_s^{(k)}) - \nabla F_k(\bar{\mathbf{x}}_s) \right), \quad B_s := \frac{1}{K} \sum_k \left( \nabla F_k(\bar{\mathbf{x}}_s) - \nabla F(\bar{\mathbf{x}}_s) \right).$$

L-smoothness gives  $\|A_s\| \leq \frac{L}{K} \sum_k \|\mathbf{x}_s^{(k)} - \bar{\mathbf{x}}_s\|$  and hence  $\mathbb{E}\|A_s\|^2 \leq L^2 D_s$ . Assumption 1 gives  $\mathbb{E}\|B_s\|^2 \leq \Gamma^2$ . Combine the three bounds.

Lemma 3 (Disagreement growth within a round) Starting from  $\mathbf{x}_0^{(k)} = \bar{\mathbf{x}}_0$ ,

$$D_s(t) \le c_0 \eta^2 s (\Gamma^2 + \sigma^2), \qquad s = 0, 1, \dots, \tau - 1,$$

for some absolute constant  $c_0 > 0$  independent of  $\eta, \tau, K, \Gamma, \sigma$ .

Sketch. Subtract the averaged update from each local update and unroll:

$$\mathbf{x}_{s+1}^{(k)} - \bar{\mathbf{x}}_{s+1} = (\mathbf{x}_s^{(k)} - \bar{\mathbf{x}}_s) - \eta \left( \nabla F_k(\mathbf{x}_s^{(k)}) - \nabla F(\bar{\mathbf{x}}_s) \right) - \eta \left( g_s^{(k)} - \nabla F_k(\mathbf{x}_s^{(k)}) \right) + \eta \zeta_s.$$

Use L-smoothness, the heterogeneity bound, and variance bounds; then take expectations to obtain linear growth in s with factor  $\eta^2(\Gamma^2 + \sigma^2)$ .

Therefore, let  $\eta \leq 1/L$  and define  $\rho := 1 - \eta \mu \in (0,1]$ . After  $\tau$  local steps and averaging,

$$\mathbb{E} \|\mathbf{x}(t+1) - \mathbf{x}^{\star}\|^{2} = \mathbb{E} \|\bar{\mathbf{x}}_{\tau}(t) - \mathbf{x}^{\star}\|^{2} \leq \rho^{\tau} \mathbb{E} \|\mathbf{x}(t) - \mathbf{x}^{\star}\|^{2} + C_{1} \frac{\eta \sigma^{2}}{\mu K} + C_{2} \left(\eta^{2} + \frac{\eta}{\mu}\right) \tau \Gamma^{2} + \mathcal{O}(\eta^{3})$$

for absolute constants  $C_1, C_2$  depending only on  $(\mu, L, c_0)$ .

### Interpretation.

- The contraction term  $\rho^{\tau}$  improves with more local steps  $\tau$ .
- Mini-batch *noise* scales as  $\frac{\eta \sigma^2}{\mu K}$ .
- The heterogeneity drift enters as a **second moment**, hence appears with  $\eta$ . Summing across  $\tau$  steps yields the characteristic  $\eta \tau \Gamma^2$  term (up to constants).

**Tuning**  $\tau$ . With homogeneous data ( $\Gamma^2 \approx 0$ ), choosing a larger  $\tau$  saves communication (the drift is negligible). With heterogeneous data (large  $\Gamma^2$ ), keep  $\tau$  modest so the  $\eta \tau \Gamma^2$  drift does not dominate.

(b) (15 points) For ECE/ORIE 5290 students: Assuming (a) holds, for homogeneous data ( $\Gamma^2 \approx 0$ ), what  $\tau$  do you recommend? For heterogeneous data (large  $\Gamma^2$ ), how would you adjust  $\tau$ ?

Homogeneous: can use larger  $\tau$  (cheap communication, minimal drift). Heterogeneous: smaller  $\tau$  to control drift; potentially adaptive  $\tau$  based on measured disagreement.

(c) (10 points) Suppose the total training time is limited. Each local iteration costs  $c_{\text{comp}}$  time units for computation, and each synchronization costs  $c_{\text{comm}}$ . Qualitatively describe how to choose  $(\eta, \tau)$  to balance runtime efficiency and convergence accuracy.

If communication is much more expensive than computation  $(c_{\mathrm{comm}}\gg c_{\mathrm{comp}})$  and the data are homogeneous, it is preferable to choose a large  $\tau$  to save time by communicating less frequently. However, if data heterogeneity is significant, a smaller  $\tau$  should be chosen to reduce drift. The learning rate  $\eta$  should start near 1/L and be reduced gradually to avoid amplifying the  $\eta^2\tau\Gamma^2$  term. A practical heuristic is to grid search over  $\tau\in\{1,2,4,8,\ldots\}$  and select the pair  $(\eta,\tau)$  that yields the best validation loss under the fixed wall-clock budget.

### Question 5: Coding - Consensus and parallel SGD (20 points)

You will implement two small simulations.

- (A) Consensus vs. Gossip (10 points) Generate N = 20 i.i.d. initial values in [0,1]. Consider:
  - Synchronous consensus  $\mathbf{x}(t+1) = \mathbf{W}\mathbf{x}(t)$  on a ring with  $\alpha = 1/2$ .
  - Randomized gossip: pick a random edge (i, j) on the ring and average the pair.

**Plot 1:** the disagreement  $V(t) = \sum_{i} (x_i(t) - \bar{x})^2$  vs. iterations for both methods (same random seed).

Plot 2: sample trajectories of two nodes to illustrate smoothing.

- (B) Local vs. Parallel SGD (10 points) Binary logistic regression on a synthetic dataset, split evenly across K = 4 workers. Compare:
  - Parallel (synchronous) mini-batch SGD with global batch size KB.
  - Local SGD with the same local batch size B and averaging period  $\tau \in \{1, 5, 20\}$ .

Plot 3: training loss  $F(\mathbf{x}) = \frac{1}{K} \sum_{k=1}^{K} F_k(\mathbf{x})$  vs. # of communication rounds. Short analysis (3–5 sentences): discuss the effect of  $\tau$  on speed/accuracy and when local SGD matches parallel SGD.

Expected outcome: The synchronous average consensus converges faster per iteration than gossip; both are linear in expectation with rates set by spectral gap. For SGD, local SGD with small  $\tau$  and homogeneous data closely matches parallel SGD in #communication rounds; large  $\tau$  saves communication but may underperform if data shards are heterogeneous.