Distributed Optimization for Machine Learning

Lecture 15 - Communication-efficient Distributed Training - Part II

Tianyi Chen

School of Electrical and Computer Engineering Cornell Tech, Cornell University

October 20, 2025





Table of Contents

Reduce the Number of Bits via Sparsification

Reduce the Number of Workers Per Round



Sparsification

Goal: Reduce num of communicated entries by making vectors sparse.

Q: What is sparse?

Quantization:



Sparsification:



Idea: Communicate only a few coordinates and set the rest to zero.



Stochastic sparsification

For any $\mathbf{v} \in \mathbb{R}^d$, define a sparsified vector $Q(\mathbf{v})$ coordinate-wise by:

$$[Q(oldsymbol{v})]_j = egin{cases} rac{v_j}{p_j}, & ext{with probability } p_j, \ 0, & ext{with probability } 1-p_j, \end{cases} j=1,\ldots,d.$$

Let $\mathbf{p}=(p_1,\ldots,p_d)$ be a predetermined probability vector belonging to a simplex $(p_j\in(0,1],\sum_{j=1}^d p_j=1)$.



Performance of stochastic sparsification

Lemma. For $\mathbf{v} \in \mathbb{R}^d$ and a sparsified vector $Q(\mathbf{v})$, it follows that

- (i) Unbiasedness: $\mathbb{E}[Q(\mathbf{v})] = \mathbf{v}$ since $\mathbb{E}[[Q(\mathbf{v})]_j] = \frac{v_j}{p_i} p_j = v_j$
- (ii) Variance bound: $\mathbb{E} \big[\| Q(\mathbf{v}) \mathbf{v} \|_2^2 \big] \le \max_j \frac{1 p_j}{p_j} \| \mathbf{v} \|_2^2$

Wang, H., Sievert, S., Liu, S., Charles, Z., Papailiopoulos, D., and Wright, S. Atomo: Communication-efficient Learning via Atomic Sparsification, *NeurIPS 2018*



D1) Threshold-based rule

For any $\mathbf{v} \in \mathbb{R}^d$, denote the sparsified vector $Q(\mathbf{v})$.

Idea: Only transmit coordinates whose magnitudes exceed τ .



D2) Memory-based threshold rule

If the algorithm transmits:

Original:
$$\mathbf{v}^{(0)}, \mathbf{v}^{(1)}, \dots, \mathbf{v}^{(K)}$$

Sparsified: $Q(\mathbf{v}^{(0)}), Q(\mathbf{v}^{(1)}), \dots, Q(\mathbf{v}^{(K)})$

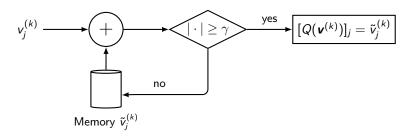
Initialize: $\tilde{\mathbf{v}}^{(0)} = \mathbf{v}^{(0)}$.

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

$$[Q(\mathbf{v}^{(k)})]_j = egin{cases} ilde{v}_j^{(k)}, & ext{if } | ilde{v}_j^{(k)}| \geq \gamma, \ 0, & ext{otherwise}. \end{cases}$$

$$\tilde{\mathbf{v}}^{(k+1)} = \mathbf{v}^{(k+1)} + (\tilde{\mathbf{v}}^{(k)} - Q(\mathbf{v}^{(k)})).$$





For
$$k = 0, 1, ..., K - 1$$
:

$$[\mathcal{Q}(\mathbf{v}^{(k)})]_j = egin{cases} ilde{v}_j^{(k)}, & ext{if } | ilde{v}_j^{(k)}| \geq \gamma, \ 0, & ext{otherwise}. \end{cases}$$

$$\tilde{\boldsymbol{v}}^{(k+1)} = \boldsymbol{v}^{(k+1)} + \left(\tilde{\boldsymbol{v}}^{(k)} - Q(\boldsymbol{v}^{(k)})\right).$$



D3) Top-k sparsification rule*

Consider $\pi \in \mathbb{R}^d$ as a permutation of $\{1,2,\ldots,d\}$ such that for $\mathbf{v} \in \mathbb{R}^d$,

$$|v_{\pi(1)}| \ge |v_{\pi(2)}| \ge \cdots \ge |v_{\pi(d)}|.$$

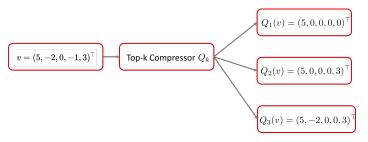
Then the j-th entry of the sparsified vector is:

$$[Q_k(\mathbf{v})]_j = egin{cases} v_j, & ext{if } j = \pi(j) ext{ and } j \leq k, \ 0, & ext{otherwise}. \end{cases}$$

Stich, S.U., Cordonnier, J.-B., and Jaggi, M. Sparsified SGD with Memory, NeurIPS 2018



D3) Top-k sparsification rule*



The error due to sparsification:

$$\|Q_k(\mathbf{v}) - \mathbf{v}\|_2^2 \le \left(1 - \frac{k}{d}\right) \|\mathbf{v}\|_2^2.$$

Stich, S.U., Cordonnier, J.-B., and Jaggi, M. Sparsified SGD with Memory, *NeurIPS* 2018

Implementation of quantized / sparsified gradient descent

For iteration $k = 1, 2, \dots, K$:

- 1. **Server broadcasts** the current model parameter x^k to all workers.
- 2. For each worker i = 1, 2, ..., n (in parallel):
 - Worker *i* calculates $\mathbf{v}_i^{(k)} = \nabla F_i(\mathbf{x}^k)$.
 - Worker *i* computes sparsified/quantized gradient $Q(\mathbf{v}_i^{(k)})$.
 - Worker *i* uploads $Q(\mathbf{v}_i^{(k)})$ to the server.
- 3. Server updates the global model:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{\alpha}{n} \sum_{i=1}^n Q(\mathbf{v}_i^{(k)}).$$

Remark: Quantization / sparsification can be performed either at the server side or at the worker side or both.



Table of Contents

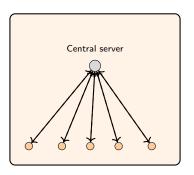
Reduce the Number of Bits via Sparsification

Reduce the Number of Workers Per Round



Reduce the number of workers

Goal: Reduce the number of workers participating in communication.

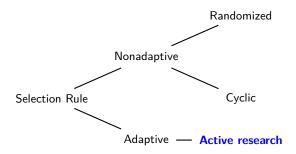


Local data on nodes



Methods for worker selection

Idea: Only a subset of workers upload/download gradients at each round, based on either fixed (nonadaptive) or dynamic (adaptive) rules.





Non-adaptive randomized rule

For iteration $k = 1, 2, \dots, K$:

- **S1)** Server randomly selects worker $i_k \in \{1, ..., n\}$ (or a set $\mathcal{I}_k \subseteq \{1, ..., n\}$), and sends \mathbf{x}^k to worker i_k (or all $i \in \mathcal{I}_k$).
 - 1. Worker i_k computes and uploads $\nabla F_{i_k}(\mathbf{x}^k)$.
 - 2. Server updates x^k via:

Option I (SGD):

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha \nabla F_{i_k}(\mathbf{x}^k).$$

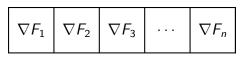
Option II (Randomized Incremental Aggregated Gradient (RIAG)):

$$\mathbf{x}_{i}^{k+1} = \begin{cases} \mathbf{x}^{k}, & i = i_{k}, \\ \mathbf{x}_{i}^{k}, & i \neq i_{k}, \end{cases} \quad \mathbf{x}^{k+1} = \mathbf{x}^{k} - \alpha \nabla F_{i_{k}}(\mathbf{x}^{k}) - \alpha \sum_{i \neq i_{k}} \nabla F_{i}(\mathbf{x}_{i}^{k}).$$



Memory overhead for RIAG

If the server pursues Option II, it stores a table $\in \mathbb{R}^{d \times n}$.



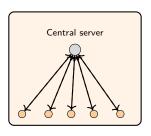
Overcome memory overhead:

Store the summation $\nabla_k = \sum_{i=1}^n \nabla F_i(\mathbf{x}_i^k)$. Worker uploads only the change of gradients:

$$\nabla_k^i = \nabla F_i(\mathbf{x}^k) - \nabla F_i(\mathbf{x}_i^k).$$

Server updates the summation via:

$$\nabla_{k+1} = \nabla_k + \nabla_k^i.$$



Local data on nodes

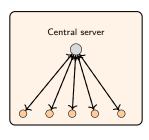


Non-adaptive cyclic rule

For k = 1, 2, ..., K:

- 1. Server selects worker $i_k = k \mod n$.
- 2. Server sends x^k to worker i_k .
- 3. Worker i_k computes and uploads $\nabla F_{i_k}(\mathbf{x}^k)$.
- 4. Server updates x^k via Option I or II.

CIAG: Cyclic Incremental Aggregated Gradient



Local data on nodes



Theoretical guarantees of CIAG

Theorem 3 (Convergence of CIAG)

Under the L-smooth and $\mu\text{-strongly}$ convex assumption, if the stepsize α in CIAG satisfies:

$$0<\alpha\leq\frac{1}{n(\mu+L)},$$

then CIAG achieves an R-linear convergence rate:

$$\|\mathbf{x}^k - \mathbf{x}^*\|_2^2 \le \rho^k \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2$$
, for some $0 < \rho < 1$.



Plan: adaptive worker selection

Compare: Gradient Descent vs. RIAG/CIAG

Tradeoff Factors:

- (c1) Amount of communication per iteration
- (c2) Number of iterations required for convergence

Observation:

RIAG/CIAG $\approx \frac{1}{n}$ communications as GD (fewer uploads per iteration), GD $\approx \frac{1}{n}$ iterations as RIAG/CIAG (faster convergence per round).

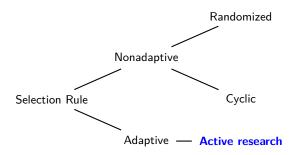
Total Communication Cost:

Total communication rounds = $(c1) \times (c2)$.



Methods for worker selection

Idea: Only a subset of workers upload/download gradients at each round, based on either fixed (nonadaptive) or dynamic (adaptive) rules.





Adaptive worker selection - best tradeoff

A slight generalization of Incremental Aggregated Gradient (IAG):

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha \sum_{i \in \mathcal{I}^k} \nabla F_i(\mathbf{x}^k) - \alpha \sum_{i \notin \mathcal{I}^k} \nabla F_i(\mathbf{x}^k_i),$$

Special cases:

- **RIAG (Randomized IAG):** $\mathcal{I}^k = \{i_k\}$, with i_k randomly generated.
- **CIAG** (Cyclic IAG): $\mathcal{I}^k = \{k \mod n\}$.
- GD (Full Gradient Descent): $\mathcal{I}^k = \{1, 2, ..., n\}$.



Incremental aggregated gradient

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha \sum_{i \in \mathcal{I}^k} \nabla F_i(\mathbf{x}^k) - \alpha \sum_{i \notin \mathcal{I}^k} \nabla F_i(\mathbf{x}^k_i)$$

$$= \underbrace{\mathbf{x}^k - \alpha \sum_{i=1}^n \nabla F_i(\mathbf{x}^k)}_{\text{GD update}} + \alpha \underbrace{\sum_{i \notin \mathcal{I}^k} \left(\nabla F_i(\mathbf{x}^k) - \nabla F_i(\mathbf{x}^k_i) \right)}_{\boldsymbol{\delta}^k_i}$$

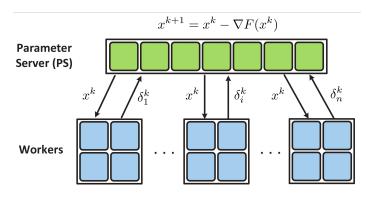
Error of using old gradients: δ_i^k

Intuition: If $\|\delta_i^k\|$ are small relative to $\sum_{i=1}^n \|\nabla F_i(\mathbf{x}^k)\|$, then the price paid for saving uploads/downloads is small.



Incremental aggregated gradient

Question: The intuition is good but how to quantify small?





Toward adaptive worker selection

Design an adaptive selection rule by analyzing the IAG iteration.

Lemma (IAG)

Under the L-smooth assumption of $F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} F_i(\mathbf{x})$, \mathbf{x}^{k+1} is generated by performing one-step generic IAG update given \mathbf{x}^k and $\{\mathbf{x}_i^k\}_{i=1}^n$, then:

$$F(\mathbf{x}^{k+1}) - F(\mathbf{x}^{k}) \leq -\frac{\alpha}{2} \|\nabla F(\mathbf{x}^{k})\|_{2}^{2} + \frac{\alpha}{2} \|\sum_{i \notin \mathcal{I}^{k}} \boldsymbol{\delta}_{i}^{k}\|_{2}^{2}$$

$$+ \left(\frac{L}{2} - \frac{1}{2\alpha}\right) \|\mathbf{x}^{k+1} - \mathbf{x}^{k}\|_{2}^{2}$$

$$\stackrel{\alpha = \frac{1}{L}}{\Longrightarrow} \leq -\frac{1}{2L} \|\nabla F(\mathbf{x}^{k})\|_{2}^{2} + \frac{1}{2L} \|\sum_{i \notin \mathcal{I}^{k}} \boldsymbol{\delta}_{i}^{k}\|_{2}^{2} \triangleq \Delta_{CIAG}^{k}.$$



Communication principle

Lemma (GD)

Under the same L-smooth assumption, the one-step GD update satisfies:

$$F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) \le -\frac{1}{2L} \|\nabla F(\mathbf{x}^k)\|_2^2 \triangleq \Delta_{GD}^k.$$

Principle: Larger progress per communication:

$$\frac{\Delta_{\mathsf{IAG}}^k}{|\mathcal{I}^k|} \leq \frac{\Delta_{\mathsf{GD}}^k}{n}$$

Plugging Δ_{GD}^k and Δ_{IAG}^k leads to:

$$\frac{-\frac{1}{2L}\|\nabla F(\boldsymbol{x}^k)\|^2 + \frac{1}{2L}\sum_{i \notin \mathcal{I}^k} \|\boldsymbol{\delta}_i^k\|^2}{|\mathcal{I}^k|} \leq \frac{-\frac{1}{2L}\|\nabla F(\boldsymbol{x}^k)\|^2}{n}$$



Deriving the sufficient condition for the principle

$$\frac{-\frac{1}{2L}\|\nabla F(\mathbf{x}^k)\|^2 + \frac{1}{2L}\sum_{i \notin \mathcal{I}^k} \|\boldsymbol{\delta}_i^k\|^2}{|\mathcal{I}^k|} \leq \frac{-\frac{1}{2L}\|\nabla F(\mathbf{x}^k)\|^2}{n}$$

$$\iff \left\|\sum_{i \notin \mathcal{I}^k} \boldsymbol{\delta}_i^k\right\|^2 \leq \left(1 - \frac{|\mathcal{I}^k|}{n}\right) \|\nabla F(\mathbf{x}^k)\|^2.$$

By Cauchy–Schwarz inequality, $\|\boldsymbol{a}_1+\boldsymbol{a}_2+\cdots+\boldsymbol{a}_n\|^2 \leq n\sum_{i=1}^n \|\boldsymbol{a}_i\|^2$, it holds that:

$$\left\| \sum_{i \notin \mathcal{I}^k} \delta_i^k \right\|^2 \le \left(n - |\mathcal{I}^k| \right) \sum_{i \notin \mathcal{I}^k} \|\delta_i^k\|^2 \le \left(n - |\mathcal{I}^k| \right) n \max_{i \notin \mathcal{I}^k} \|\delta_i^k\|^2.$$



Deriving the sufficient condition for progress principle

Sufficient Condition for the Principle:

$$(n - |\mathcal{I}^k|) n \max_{i \notin \mathcal{I}^k} \|\boldsymbol{\delta}_i^k\|^2 \le \frac{n - |\mathcal{I}^k|}{n} \|\nabla F(\boldsymbol{x}^k)\|^2.$$

$$\iff \|\boldsymbol{\delta}_i^k\|^2 \le \frac{1}{\alpha^2 n^2} \|\nabla F(\boldsymbol{x}^k)\|^2, \quad \text{for all } i \in \{1, \dots, n\}.$$

Q: How can we check this condition either at the server or at worker?

$$\|\nabla F(\mathbf{x}^k)\|^2 = \left\|\sum_{i=1}^n \nabla F_i(\mathbf{x}^k)\right\|^2$$

This cannot be computed locally.



Checking the sufficient condition

Approximation:

$$\|\nabla F(\mathbf{x}^k)\|^2 \approx \frac{1}{\alpha^2} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2$$

so that each worker can check condition locally by:

$$\|\boldsymbol{\delta}_i^k\|^2 \leq \frac{1}{lpha^2 n^2} \|\boldsymbol{x}^k - \boldsymbol{x}^{k-1}\|^2$$
 (Worker side)

Q: What if we find an upper bound on the left-hand side?

$$\|\nabla F_i(\mathbf{x}^k) - \nabla F_i(\mathbf{x}_i^k)\| \le L_i \|\mathbf{x}^k - \mathbf{x}_i^k\|$$

A sufficient condition rule is:

$$||L_i^2||\mathbf{x}^k - \mathbf{x}_i^k||^2 \le \frac{1}{\alpha^2 n^2} ||\mathbf{x}^k - \mathbf{x}^{k-1}||^2$$
 (Server side)



Implementation of adaptive selection rule (LAG)*

Worker side:

For iteration $k = 1, 2, \dots, K$:

- 1. Server broadcasts the current model parameter x^k to all workers.
- 2. For each worker i = 1, 2, ..., n (in parallel):
 - Worker *i* computes the local gradient $\nabla F_i(\mathbf{x}^k)$.
 - Worker i checks the upload condition:

$$\|\boldsymbol{\delta}_{i}^{k}\|^{2} \leq \frac{1}{\alpha^{2}n^{2}}\|\boldsymbol{x}^{k} - \boldsymbol{x}^{k-1}\|^{2}.$$

- If the condition is satisfied ⇒ Do not upload.
- Otherwise ⇒ Upload.
- 3. **Server updates** the global model via the generic IAG update rule.

Chen, T., Giannakis, G., Sun, T., and Yin, W. LAG: Lazily Aggregated Gradient for Communication-efficient Distributed Learning, *NeurIPS 2018*



Implementation of adaptive selection rule (LAG)*

Server side:

For iteration $k = 1, 2, \dots, K$:

1. **Server checks** the condition for each worker i = 1, 2, ..., n:

$$||\mathbf{L}_{i}^{2}||\mathbf{x}^{k}-\mathbf{x}_{i}^{k}||^{2} \leq \frac{1}{\alpha^{2}n^{2}}||\mathbf{x}^{k}-\mathbf{x}^{k-1}||^{2}.$$

- 2. Collect all violating workers into the set \mathcal{I}^k .
- 3. **Server sends** the current model x^k to all $i \in \mathcal{I}^k$.
- 4. For each worker $i \in \mathcal{I}^k$:
 - Worker *i* computes and uploads $\nabla F_i(\mathbf{x}^k)$ to the server.
- 5. **Server updates** the global parameter x^{k+1} via the generic IAG update rule.

Chen, T., Giannakis, G., Sun, T., and Yin, W. LAG: Lazily Aggregated Gradient for Communication-efficient Distributed Learning, *NeurIPS 2018*



Theoretical guarantee of LAG

Theorem 4 (Convergence of LAG)

1. Under the L-smooth assumption of $F_i(x)$, we have:

$$\frac{1}{K} \sum_{k=1}^{K} \|\nabla F(\mathbf{x}^k)\|^2 = \mathcal{O}\left(\frac{1}{K}\right)$$
 (Same as GD)

2. Under the additional convex assumption, we have:

$$F(\mathbf{x}^k) - F(\mathbf{x}^*) = \mathcal{O}\left(\frac{1}{K}\right)$$
 (Same as GD)

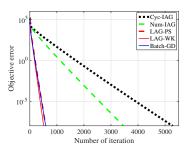
3. Under the additional μ -strong convexity assumption, we have:

$$F(\mathbf{x}^k) - F(\mathbf{x}^*) = \mathcal{O}((1 - \frac{\mu}{L})^k)$$
 (Same as GD)



Empirical performance of LAG

- Faster convergence per iteration: LAG achieves similar or faster convergence compared with IAG and GD in terms of iteration complexity.
- Significantly reduced communication cost: LAG requires fewer communication rounds while maintaining accuracy.



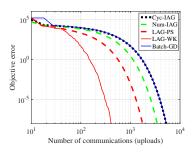




Figure: Iteration and communication complexity for linear regression.

Worker-side condition in LAG (setup)

Goal: Show that Lazy Aggregated Gradient (LAG) still guarantees descent even when workers communicate intermittently.

Worker update condition:

$$\|\boldsymbol{\delta}_i^k\|^2 \leq \frac{\zeta}{\alpha^2 n^2} \|\boldsymbol{x}^k - \boldsymbol{x}^{k-1}\|^2,$$

where $\boldsymbol{\delta}_{i}^{k} = \nabla F_{i}(\boldsymbol{x}_{i}^{k}) - \nabla F_{i}(\boldsymbol{x}^{k})$.

Intuition:

- lacksquare δ_i^k measures how "stale" a worker's gradient is.
- Larger $\zeta \Rightarrow$ easier condition \Rightarrow fewer communications.
- When local changes are small, workers can skip communication.



Step 1: One-step progress of LAG

Under L-smoothness of F, the descent lemma gives:

$$F(\boldsymbol{x}^{k+1}) - F(\boldsymbol{x}^k) \leq \langle \nabla F(\boldsymbol{x}^k), \boldsymbol{x}^{k+1} - \boldsymbol{x}^k \rangle + \frac{L}{2} \|\boldsymbol{x}^{k+1} - \boldsymbol{x}^k\|^2.$$

Substitute the LAG update rule:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha \sum_{i \in \mathcal{T}^k} \nabla F_i(\mathbf{x}^k) - \alpha \sum_{i \notin \mathcal{T}^k} \nabla F_i(\mathbf{x}^k_i).$$

This introduces gradient mismatch terms δ_i^k from lazy workers.



Bounding the inner product term (sketch)

Plugging the update into the inner product:

$$\langle \nabla F(\mathbf{x}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle = -\alpha \|\nabla F(\mathbf{x}^k)\|^2 - \alpha \langle \nabla F(\mathbf{x}^k), \sum_{i \notin \mathcal{I}^k} \delta_i^k \rangle.$$

Using the identity $2a^{T}b = ||a||^{2} + ||b||^{2} - ||a - b||^{2}$, we bound the cross term:

$$\left\langle \nabla F(\mathbf{x}^k), \sum_{i \notin \mathcal{I}^k} \delta_i^k \right\rangle \leq \frac{1}{2} \|\nabla F(\mathbf{x}^k)\|^2 + \frac{1}{2} \|\sum_{i \notin \mathcal{I}^k} \delta_i^k \|^2.$$

Result:

$$F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) \le -\frac{\alpha}{2} \|\nabla F(\mathbf{x}^k)\|^2 + \frac{\alpha}{2} \sum_{i \neq \tau k} \|\boldsymbol{\delta}_i^k\|^2 + \frac{L}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2.$$



Step 2: Apply worker-side condition

Applying the worker condition to bound the error term gives:

$$F(\mathbf{x}^{k+1}) - F(\mathbf{x}^k) \le -\frac{\alpha}{2} \|\nabla F(\mathbf{x}^k)\|^2 + \frac{\zeta}{2\alpha n} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2.$$

Takeaway:

- The first term guarantees descent.
- The second term captures accumulated gradient staleness.
- lacktriangle Properly tuning ζ keeps LAG stable and communication-efficient.



Understanding the tradeoff

Key tradeoff:

Descent rate vs. Error from delayed gradients.

If ζ is too small \Rightarrow frequent communication but stable. If ζ is too large \Rightarrow fewer communications but possible instability.

Balance point: choose $\zeta = 1 - \alpha L$ so that:

Error term \sim Descent term.

This ensures that each step still decreases F(x) in expectation.



Step 3: Convergence sketch

Choose: $\zeta = 1 - \alpha L$ to balance descent and error.

Telescoping over k = 1, ..., K:

$$\frac{1}{K}\sum_{k=1}^K \|\nabla F(\mathbf{x}^k)\|^2 \leq \frac{F(\mathbf{x}^1) - F(\mathbf{x}^*)}{\alpha K}.$$

Result:

LAG achieves
$$\mathcal{O}(1/\mathcal{K})$$
 convergence

- the same as parallel GD, but with significantly fewer communications.

(Idea: smoothness \Rightarrow gradients evolve slowly \Rightarrow lazy communication.)



Recap and fine-tuning

- What we have talked about today?
 - \Rightarrow **Sparsification** reduces communication cost by transmitting gradients with fewer entries.
 - ⇒ Worker selection reduces communication cost by letting only a subset of workers upload gradients adaptively.



Welcome anonymous survey!



