

Stochastic Extragradient with Flip-Flop Anchoring: Provable Improvements

Jiseok Chae¹ Chulhee Yun² Donghwan Kim¹

¹Department of Mathematical Sciences, KAIST

²Kim Jaechul Graduate School of AI, KAIST

NeurIPS 2024

Minimax Problems

We consider unconstrained minimax problems with a *finite-sum* structure:

$$\min_{\mathbf{x}} \max_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}, \mathbf{y}).$$

Very versatile, and has many ML applications:

- Generative Adversarial Networks
- Consistency Trajectory Models
- Sharpness-aware Minimization
- Computing Optimal Transport Maps
- ...

Minimax Problems

We consider unconstrained minimax problems with a *finite-sum* structure:

$$\min_{\mathbf{x}} \max_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}, \mathbf{y}).$$

Denote both min. and max. variables at once by $\mathbf{z} := (\mathbf{x}, \mathbf{y})$.

The *saddle gradient*

$$\mathbf{F}(\mathbf{z}) = \begin{bmatrix} \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) \\ -\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) \end{bmatrix}$$

is more natural than ∇f in minimax problems.

The Extragradient Method

The *gradient descent-ascent* (GDA) method

$$\mathbf{z}_{k+1} = \mathbf{z}_k - \eta_k \mathbf{F}(\mathbf{z}_k) \quad \text{or} \quad \begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \eta_k \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k) \\ \mathbf{y}_{k+1} &= \mathbf{y}_k + \eta_k \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k) \end{aligned}$$

already does not work for simple convex-concave problems.

The *extragradient* (EG) method (Korpelevich, 1976)

$$\begin{aligned} \mathbf{w}_k &= \mathbf{z}_k - \eta_k \mathbf{F}(\mathbf{z}_k) \\ \mathbf{z}_{k+1} &= \mathbf{z}_k - \eta_k \mathbf{F}(\mathbf{w}_k) \end{aligned} \quad \text{or} \quad \mathbf{z}_{k+1} = \mathbf{z}_k - \eta_k \mathbf{F}(\mathbf{z}_k - \eta_k \mathbf{F}(\mathbf{z}_k))$$

on the other hand, **works** on convex-concave problems.

Stochastic Extragradient?

Unlike GDA vs EG, the *stochastic* EG (SEG)

$$\mathbf{z}_{k+1} = \mathbf{z}_k - \eta_k \mathbf{F}_{i(k)}(\mathbf{z}_k - \eta_k \mathbf{F}_{i(k)}(\mathbf{z}_k))$$

does not show a clear advantage in convex-concave problems over GDA.

Even if we additionally assume each f_i are also convex-concave, convergence rates typically look something like:

$$\min_{k=0,1,\dots,K} \|\mathbf{F}\mathbf{z}_k\|^2 \leq \mathcal{O}\left(\frac{1}{\text{poly}(K)}\right) + (\text{abs. const.})$$

* The constant term can be decreased only with strong additional assumptions, such as for example, increasing the batch size every iteration.

For minimization problems...

With-replacement stochastic gradient descent (SGD) works well.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \nabla f_{i(k)}(\mathbf{x}_k), \quad i(k) \sim \text{Unif}(\{1, \dots, n\})$$

In practice, *shuffling based* SGD is used.

Random reshuffling (RR): in the k th epoch, a permutation $\tau_k : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ is chosen randomly, and

$$\begin{aligned} \mathbf{x}_i^k &= \mathbf{x}_{i-1}^k - \eta_k \nabla f_{\tau_k(i)}(\mathbf{x}_{i-1}^k), \quad i = 1, \dots, n, \\ \mathbf{x}_0^{k+1} &\leftarrow \mathbf{x}_n^k. \end{aligned}$$

For minimization problems...

With-replacement stochastic gradient descent (SGD) works well.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \nabla f_{i(k)}(\mathbf{x}_k), \quad i(k) \sim \text{Unif}(\{1, \dots, n\})$$

In practice, *shuffling based* SGD is used.

Flip-flop sampling (FF) (Rajput et al., 2022) goes one step even further in search for a better sampling scheme: in the k th epoch, a permutation $\tau_k : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ is chosen randomly, and

$$\begin{aligned} \mathbf{x}_i^k &= \mathbf{x}_{i-1}^k - \eta_k \nabla f_{\tau_k(i)}(\mathbf{x}_{i-1}^k), & i &= 1, \dots, n, \\ \mathbf{x}_i^k &= \mathbf{x}_{i-1}^k - \eta_k \nabla f_{\tau_k(2n+1-i)}(\mathbf{x}_{i-1}^k), & i &= n+1, \dots, 2n, \\ \mathbf{x}_0^{k+1} &\leftarrow \mathbf{x}_{2n}^k. \end{aligned}$$

For minimization problems...

With-replacement stochastic gradient descent (SGD) works well.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \nabla f_{i(k)}(\mathbf{x}_k), \quad i(k) \sim \text{Unif}(\{1, \dots, n\})$$

In practice, *shuffling based* SGD is used.

In terms of convergence rates,

- RR is in general faster than with-replacement SGD.

(Ahn et al., 2020; Mishchenko et al., 2020)

- If all f_i are quadratic functions then FF is even faster, thanks to the stochastic error term being smaller.

(Rajput et al., 2022)

Our contributions

- Stochastic EG indeed does not work on convex-concave problems. Shuffling does not resolve the problem.
 - An explicit counterexample with divergent iterates
- On top of FF, adding a simple *anchoring* step

$$z_0^{k+1} \leftarrow \frac{z_{2n}^k + z_0^k}{2}$$

reduces the stochastic error by an order of magnitude (w.r.t. stepsize), finally allowing a convergence rate of $\tilde{O}(1/k^{1/3})$.

- The reduced error also benefits the convergence on strongly-convex-strongly-concave problems, enjoying a rate of $\tilde{O}(1/nk^4)$.
 - Without anchoring (*i.e.*, with-replacement sampling or RR only), the convergence rate is at best $\Omega(1/nk^3)$.

Algorithm

Stochastic Extragradient with Flip-Flop Anchoring (SEG-FFA)

For each $k = 0, 1, \dots$: # epoch level outer loop

$\tau_k \sim \text{Unif}(\mathfrak{S}_n)$ # sample random permutation

For each $i = 1, \dots, n$: # flip

$$z_i^k = z_{i-1}^k - \eta_k \mathbf{F}_{\tau_k(i)} \left(z_{i-1}^k - \frac{\eta_k}{2} \mathbf{F}_{\tau_k(i)}(z_{i-1}^k) \right)$$

For each $i = n+1, \dots, 2n$: # flop

$$z_i^k = z_{i-1}^k - \eta_k \mathbf{F}_{\tau_k(2n+1-i)} \left(z_{i-1}^k - \frac{\eta_k}{2} \mathbf{F}_{\tau_k(2n+1-i)}(z_{i-1}^k) \right)$$

$z_0^{k+1} \leftarrow \frac{z_{2n}^k + z_0^k}{2}$ # anchoring

Thank you for your attention.

Visit us at the Poster Session!

References I

- Ahn, Kwangjun, Chulhee Yun, and Suvrit Sra (2020). “SGD with shuffling: optimal rates without component convexity and large epoch requirements”. In: *Advances in Neural Information Processing Systems* 33, pp. 17526–17535.
- Korpelevich, Galina M (1976). “The extragradient method for finding saddle points and other problems”. In: *Matecon* 12, pp. 747–756.
- Mishchenko, Konstantin, Ahmed Khaled, and Peter Richtárik (2020). “Random reshuffling: Simple analysis with vast improvements”. In: *Advances in Neural Information Processing Systems* 33, pp. 17309–17320.
- Rajput, Shashank, Kangwook Lee, and Dimitris Papailiopoulos (2022). “Permutation-Based SGD: Is Random Optimal?” In: *International Conference on Learning Representations*.