

Hessian-Amended Random Perturbation (HARP) Using Noisy Zeroth-Order Oracle

Jingyi Zhu

MSML21: Mathematical and Scientific Machine Learning
Session 1: Optimization and Algorithms

August 16, 2021

Minimization Using *Few* Zeroth-Order Queries

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} L(\boldsymbol{\theta}) \equiv \mathbb{E}_{\boldsymbol{\omega} \sim \mathbb{P}} [\ell(\boldsymbol{\theta}, \boldsymbol{\omega})], \quad (1)$$

- **stochastic**: evaluation of $L(\boldsymbol{\theta})$ is corrupted by *noise*
- **limited-resource**: collecting $\ell(\cdot, \boldsymbol{\omega})$ is *expensive*

Stochastic Approximation (SA) Algorithms

$$\text{1st-order : } \hat{\boldsymbol{\theta}}_{k+1} = \hat{\boldsymbol{\theta}}_k - a_k \hat{\mathbf{g}}_k(\hat{\boldsymbol{\theta}}_k). \quad (2)$$

w/ a_k is stepsize. We use gradient estimator using two ZO queries:

$$\hat{\mathbf{g}}_k(\hat{\boldsymbol{\theta}}_k) = \frac{\ell(\hat{\boldsymbol{\theta}}_k + c_k \boldsymbol{\Delta}_k, \boldsymbol{\omega}_k^+) - \ell(\hat{\boldsymbol{\theta}}_k - c_k \boldsymbol{\Delta}_k, \boldsymbol{\omega}_k^-)}{2c_k} \mathbf{m}_k(\boldsymbol{\Delta}_k). \quad (3)$$

w/ c_k is differencing magnitude, $\boldsymbol{\Delta}_k$ is perturbation, $\mathbf{m}_k(\cdot) : \mathbb{R}^d \mapsto \mathbb{R}^d$.

RDSA (random direction) [Erm69, Erm83]

$$\Delta_k \sim \text{Unif}(\mathbb{S}) \quad \text{and} \quad \mathbf{m}_k(\Delta_k) = d\Delta_k$$

SPSA (simultaneous perturbation) [Spa92]

$$\Delta_k \sim [\text{Unif}\{-1, 1\}]^d, \quad \mathbf{m}_k(\Delta_k) = \Delta_k$$

SFSA (smoothed functional) [KO72]

$$\Delta_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad \text{and} \quad \mathbf{m}_k(\Delta_k) = \Delta_k$$

Randomized FDSA (finite-difference)—effective as cyclic scheme

$$\Delta_k \sim \text{Unif}\{\mathbf{e}_1, \dots, \mathbf{e}_d\} \quad \text{and} \quad \mathbf{m}_k(\Delta_k) = d\Delta_k$$

What if Δ_k has zero-mean and Σ_k^{-1} -covariance for $\Sigma_k \succ \mathbf{0}$?

- Data-driven Σ_k can be \mathcal{F}_k -measurable random variable
- User-specified Σ_k can use domain knowledge

Scaling/Stepsize

Minimize $L(\theta) = 100\theta_1^2 + \theta_2^2$, with $\hat{\theta} = [1, 1]^T$ and $c = 0.1$.

- SPSA generates $[1, 1]^T$, $[1, -1]^T$, $[-1, 1]^T$ and $[-1, -1]^T$ equally likely. $\mathbb{E}_{\Delta}[\hat{g}(\hat{\theta})]$ equals true gradient $g(\hat{\theta}) = [100, 1]^T$. But the variance is in the order of 10^4 assuming noise-free ZO queries.
- HARP draws Δ from $\mathbf{0}$ -mean and Σ^{-1} -covariance distribution with $\Sigma = \hat{H}(\hat{\theta})$. Still unbiased, but covariance matrix norm is 2×10^2 .

Correlation/Direction

- $\Sigma = \mathbf{I}$ has zero off-diagonal elements implies that each component of Δ is *independent*.
- Say $L(\theta) = 100\theta_1^2 + \theta_2^2 + \theta_1\theta_2$ with same $\hat{\theta}$ and c . The covariance magnitude of SPSA gradient estimator is around 4×10^4 and that of HARP is 8×10^2 .

HARP handles scaling & correlation

Δ_k follows a dist. w/ mean $\mathbf{0}$ and cov. $\hat{\mathbf{H}}_k^{-1}$, and $\mathbf{m}_k(\Delta_k) = \hat{\mathbf{H}}_k \Delta_k$

Theoretical Guarantee

Root-mean-squared error $\mathbb{E}[\|\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*\|^2]^{1/2}$ is **smaller** when $\boldsymbol{\Sigma}_k = \hat{\mathbf{H}}_k$ than when $\boldsymbol{\Sigma}_k = \mathbf{I}$ for **ill-conditioned** problem.

- Expectation of estimation error goes to zero at the same^a rate.
- Variance gets smaller while using $\boldsymbol{\Sigma}_k = \hat{\mathbf{H}}_k$.

^a $k^{-1/3}$ when ω_k^+ and ω_k^- are independent and identically distributed (IID).
 $k^{-1/2}$ when $\omega_k^+ = \omega_k^-$, referred to as common random number (CRN).

Hessian can be estimated [Spa00, Spa09] provides principled way to estimate Hessian using *four* loss function evaluations. [Zhu21] proposes other forms that uses *two or more*.

However, computing issue persists, though [ZWS19] reduces per-iteration FLOPs from $O(d^3)$ to $O(d^2)$. Not comparable with $O(d)$ for generic first-order methods.

Feed Hessian estimate into both Σ_k

- ① **generate** Δ_k so that it has a mean of $\mathbf{0}$ and a covariance \hat{H}_k^{-1}
- ② **collect** two noisy losses and **estimate** \hat{g}_k using
$$m_k(\Delta_k) \left[\ell(\hat{\theta}_k + c_k \Delta_k) - \ell(\hat{\theta}_k - c_k \Delta_k) \right] / (2c_k)$$
- ③ *may* **collect** additional queries and **estimate** \hat{H}_k

Skew-quartic function is ill-conditioned: one large eigenvalues, and remaining eigenvalues are close to zero.

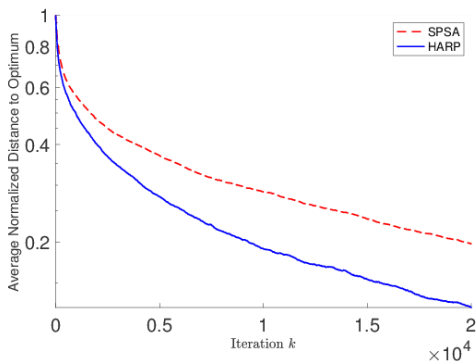


Figure: Performance of SPSA and HARP in terms of normalized distance $\|\hat{\theta}_k - \theta^*\| / \|\hat{\theta}_0 - \theta^*\|$ averaged across 25 independent replicates, and both algorithms use four ZO queries per iteration. The underlying loss function is the skew-quartic function with $d = 20$, and the noisy observation is corrupted by a $\mathcal{N}(0, 1)$ noise.

We consider generating adversarial perturbation universally for $I > 1$ images [CZS⁺17, CLC⁺18]:

$$\left\{ \begin{array}{l} \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \underbrace{\kappa \|\boldsymbol{\theta}\|_2^2}_{\equiv L_1(\boldsymbol{\theta})} + \underbrace{\frac{1}{I} \sum_{i=1}^I \text{loss}(\zeta_i + \boldsymbol{\theta})}_{\equiv L_2(\boldsymbol{\theta})}, \\ \text{s.t. } (\zeta_i + \boldsymbol{\theta}) \in [-0.5, 0.5]^d, \forall i, \end{array} \right. \quad (4)$$

where the constraint is to normalize the resulting pixels within $[-0.5, 0.5]^d$, and $\text{loss}(\cdot) : \mathbb{R}^d \mapsto \mathbb{R}$ on each image in (4) follows from [CW17]. Note that $L_2(\boldsymbol{\theta}) = 0$ when all the selected images are successfully attacked. The noisy loss observation $\ell(\boldsymbol{\theta}, \omega)$ in (1) is

$$\ell(\boldsymbol{\theta}, \omega) = \kappa \|\boldsymbol{\theta}\|_2^2 + \frac{1}{J} \sum_{j=1}^J \text{loss}(\zeta_{i_j(\omega)} + \boldsymbol{\theta}), \quad (5)$$

for $J \leq I$, and the J indexes $\{i_1(\omega), \dots, i_J(\omega)\}$ are i.i.d. uniformly drawn from $\{1, \dots, I\}$ (without replacement).

Algo	$\mathbb{E}[L(\hat{\theta}_K)]$	$\{\text{Var}[L(\hat{\theta}_K)]\}^{\frac{1}{2}}$	$\mathbb{E}[L_2(\hat{\theta}_K)]$
ADAMM	185.96	16.88	40.95
HARP	138.22	18	12.50

Table: Performance of ZO-ADAMM and HARP in terms of loss after $K = 1000$ iterations averaged across 25 independent replicates. The loss function $L(\cdot)$ is the sum of the magnitude cost $L_1(\cdot)$ and the attack loss $L_2(\cdot)$. Here $L_2(\cdot)$ measures the attack loss on $I = 100$ images of the letter one, and its *noisy* query is evaluated using a batch-size of one. A close-to-zero $L_2(\cdot)$ loss is equivalent to a close-to-one attack success rate.

Algo	$\mathbb{E}[L(\hat{\theta}_K)]$	$\{\text{Var}[L(\hat{\theta}_K)]\}^{\frac{1}{2}}$	$\mathbb{E}[L_2(\hat{\theta}_K)]$
ADAMM	56.95	6.89	11.75
HARP	18.46	1.37	0.13

Table: Here $L_2(\cdot)$ measures the attack loss on $I = 10$ images of the letter three, and its ZO query is noise-free.


Summary

We propose HARP that feeds second-order approximation to Σ_k , less sensitive to ill-conditioning.


- framework allowing \mathcal{F}_k -measurable perturbation covariance
- second-order info not only gets into parameter update but also *search scaling/direction* in HARP
- asym. rate of convergence remains the same, yet RMS is smaller

Future Work

- this framework can be generalized to scenario where 1st-order oracles are available—Hessian can be estimated using three noisy gradients
- other user-specified structure for Σ_k




Minhao Cheng, Thong Le, Pin-Yu Chen, Huan Zhang, JinFeng Yi, and Cho-Jui Hsieh.
Query-efficient hard-label black-box attack: An optimization-based approach.
In *International Conference on Learning Representations*, 2018.




Nicholas Carlini and David Wagner.
Towards evaluating the robustness of neural networks.
In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017.



Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh.
Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models.
In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, pages 15–26, 2017.




Yu M Ermol'ev.
On the method of generalized stochastic gradients and quasi-fejér sequences.
Cybernetics, 5(2):208–220, 1969.




Yuri Ermoliev.
Stochastic quasigradient methods and their application to system optimization.
Stochastics: An International Journal of Probability and Stochastic Processes, 9(1-2):1–36, 1983.



V Ya Katkovnik and KULCHITS. OY.
Convergence of a class of random search algorithms.
Automation and Remote Control, 33(8):1321–1326, 1972.



James C. Spall.
Multivariate stochastic approximation using a simultaneous perturbation gradient approximation.
IEEE transactions on automatic control, 37(3):332–341, 1992.



James C Spall.
Adaptive stochastic approximation by the simultaneous perturbation method.
IEEE transactions on automatic control, 45(10):1839–1853, 2000.



James C Spall.

Feedback and weighting mechanisms for improving jacobian estimates in the adaptive simultaneous perturbation algorithm.

IEEE Transactions on Automatic Control, 54(6):1216–1229, 2009.



Jingyi Zhu.

Hessian estimation via stein's identity in black-box problems.

In 2nd Conference on Mathematical and Scientific Machine Learning, 2021.



J. Zhu, L. Wang, and J. C. Spall.

Efficient implementation of second-order stochastic approximation algorithms in high-dimensional problems.

IEEE Transactions on Neural Networks and Learning Systems, in press at <http://dx.doi.org/10.1109/TNNLS.2019.2935455>, 2019.