

Stochastic Newton methods with enhanced Hessian estimation

Danda Sai Koti Reddy

M.Sc(engg) thesis defence

Research Advisor : Prof. Shalabh Bhatnagar

Stochastic Systems Lab,
Department of Computer Science and Automation,
Indian Institute of Science, Bangalore

Simulation Optimization

Random Directions Stochastic Approximation (RDSA) +
Improved Hessian Estimation

Improved Hessian Estimation for Simultaneous Perturbation
Stochastic Approximation with 3 Simulations (SPSA-3)

Numerical Results

Simulation Optimization

Energy Demand Management

- Consumer demand, energy generation are uncertain
- Objective is to minimize absolute difference



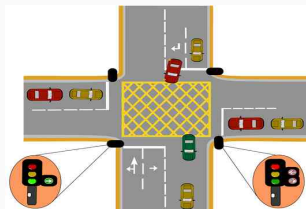
Energy Demand Management

- Consumer demand, energy generation are uncertain
- Objective is to minimize absolute difference



Traffic Signal Control

- Optimal order to switch traffic lights
- Objective is to minimize waiting time



Basic optimization problem

To find θ^* that minimizes the objective function $f(\theta)$:

$$\theta^* = \arg \min_{\theta \in \Theta} f(\theta) \quad (1)$$

- $f: \mathbb{R}^N \rightarrow \mathbb{R}$ is called the **objective function**
- θ is tunable N-dimensional parameter
- $\Theta \subseteq \mathbb{R}^N$ is the **feasible region** in which θ takes values

Deterministic optimization problem

- Complete information about objective function f
- First and higher order derivatives
- Feasible region

Classification of optimization problems

Deterministic optimization problem

- Complete information about objective function f
- First and higher order derivatives
- Feasible region

Stochastic optimization problem

- We have little knowledge on the structure of f
- f cannot be obtained directly
- $f(\theta) \equiv E_{\xi}[h(\theta, \xi)]$, where ξ comprises the randomness in the system

Difficult to find θ^* only on the basis of noisy samples

Stochastic optimization via simulation

Stochastic optimization deals with highly nonlinear and high dimensional systems. The challenges with these problems are:

- Too complex to solve analytically
- Many simplifying assumptions are required

Stochastic optimization via simulation

Stochastic optimization deals with highly nonlinear and high dimensional systems. The challenges with these problems are:

- Too complex to solve analytically
- Many simplifying assumptions are required

A good alternative of modelling and analysis is “Simulation”

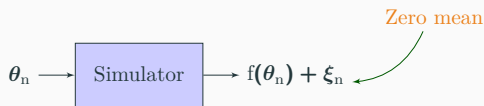


Figure 1: Simulation optimization

Stochastic analog of gradient descent

$$\theta_{n+1} = \Gamma_{\Theta} \left[\theta_n - a_n \widehat{\nabla} f(\theta_n) \right] \quad (2)$$

- $\widehat{\nabla} f(\theta_n)$ is a **noisy** estimate of the gradient $\nabla f(\theta_n)$, and it should satisfy $E \left[\widehat{\nabla} f(\theta_n) \right] - \nabla f(\theta_n) \rightarrow 0$

- $\{a_n\}$ are **pre-determined** step-sizes satisfying:

$$\sum_{n=1}^{\infty} a_n = \infty, \quad \sum_{n=1}^{\infty} a_n^2 < \infty$$

- Γ_{Θ} denotes the projection of a point onto Θ

Related second-order methods

(Spall 2000) ¹	Second-order SPSA (2SPSA)	4 simulations/iteration
(Spall 2009) ²	2SPSA + feedback	4 simulations/iteration
(Prashanth L.A. et al 2016) ³	Second-order RDSA (2RDSA)	3 simulations/iteration
(S. Bhatnagar et al 2015) ⁴	Second-order SPSA-3 (2SPSA-3)	3 simulations/iteration

¹J. C. Spall (2000), “Adaptive stochastic approximation by the simultaneous perturbation method,” IEEE TAC.

²J. C. Spall (2009), “Feedback and weighting mechanisms for improving Jacobian estimates in the adaptive simultaneous perturbation algorithm,” IEEE TAC.

³Prashanth L. A, Shalabh Bhatnagar, Michael Fu, Steve Marcus (2016), “Adaptive system optimization using random directions stochastic approximation,” IEEE TAC.

⁴S. Bhatnagar, Prashanth L. A (2015) ,“Simultaneous perturbation Newton algorithms for simulation optimization,” JOTA.

Our work

- We propose generalised RDSA algorithm¹ + feedback and weighting mechanisms for improving Hessian estimate²
- We propose feedback and weighting mechanisms for improving Hessian estimate of 2SPSA-3 algorithm¹

¹Under preparation - “<https://github.com/dsai1215/asgoodasitgets/tree/master/Journal>”.

²D. Sai Koti Reddy, Prashanth L.A, Shalabh Bhatnagar (2016), “Improved Hessian estimation for adaptive random directions stochastic approximation,” IEEE CDC 2016.

Random Directions Stochastic
Approximation (RDSA) + Improved
Hessian Estimation

Our algorithm

- Matrix projection
- Gradient estimate

$$\theta_{n+1} = \theta_n - a_n \gamma (\bar{H}_n)^{-1} \hat{\nabla} f(\theta_n) \quad (3)$$


Our algorithm

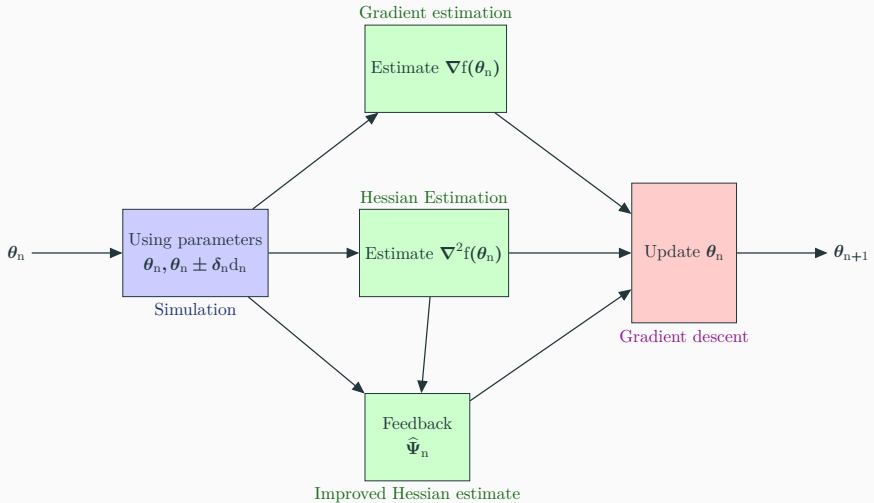
- Matrix projection
- Gradient estimate

$$\theta_{n+1} = \theta_n - a_n \gamma (\bar{H}_n)^{-1} \hat{\nabla} f(\theta_n) \quad (3)$$

$$\bar{H}_n = (1 - b_n) \bar{H}_{n-1} + b_n (\hat{H}_n - \hat{\Psi}_n) \quad (4)$$

- Optimal step-sizes
- Hessian estimate
- Feedback term

Overall flow of 2RDSA-IH



Function measurements

$$y_n^+ = f(\theta_n + \delta_n d_n) + \xi_n^+, \quad y_n^- = f(\theta_n - \delta_n d_n) + \xi_n^-$$

Function measurements

$$y_n^+ = f(\theta_n + \delta_n d_n) + \xi_n^+, \quad y_n^- = f(\theta_n - \delta_n d_n) + \xi_n^-$$

Gradient estimate

$$\hat{\nabla}f(\theta_n) = \frac{1}{\lambda} d_n \left[\frac{y_n^+ - y_n^-}{2\delta_n} \right] \quad (5)$$

Where $d_n = (d_n^1, \dots, d_n^N)^T$, and $\lambda = \mathbb{E}(d_n^i)^2$

Function measurements

$$y_n^+ = f(\theta_n + \delta_n d_n) + \xi_n^+, \quad y_n^- = f(\theta_n - \delta_n d_n) + \xi_n^-, \quad y_n = f(\theta_n) + \xi_n$$

Function measurements

$$y_n^+ = f(\theta_n + \delta_n d_n) + \xi_n^+, \quad y_n^- = f(\theta_n - \delta_n d_n) + \xi_n^-, \quad y_n = f(\theta_n) + \xi_n$$

Hessian estimate \hat{H}_n

$$\begin{aligned} \hat{H}_n &= M_n \left(\frac{y_n^+ + y_n^- - 2y_n}{\delta_n^2} \right) \\ &= M_n \left[\left(\frac{f(\theta_n + \delta_n d_n) + f(\theta_n - \delta_n d_n) - 2f(\theta_n)}{\delta_n^2} \right) \right. \\ &\quad \left. + \left(\frac{\xi_n^+ + \xi_n^- - 2\xi_n}{\delta_n^2} \right) \right] \\ &= M_n \left(d_n^T \nabla^2 f(\theta_n) d_n + O(\delta_n^2) + \left(\frac{\xi_n^+ + \xi_n^- - 2\xi_n}{\delta_n^2} \right) \right) \quad (6) \end{aligned}$$

Want to recover

 $\nabla^2 f(\theta_n)$ from this

Zero-mean

How to choose M_n ?

$$M_n = \begin{bmatrix} \frac{1}{\kappa} ((d_n^1)^2 - \lambda) & \cdots & \frac{1}{2\lambda^2} d_n^1 d_n^N \\ \frac{1}{2\lambda^2} d_n^2 d_n^1 & \cdots & \frac{1}{2\lambda^2} d_n^2 d_n^N \\ \cdots & \cdots & \cdots \\ \frac{1}{2\lambda^2} d_n^N d_n^1 & \cdots & \frac{1}{\kappa} ((d_n^N)^2 - \lambda) \end{bmatrix} \quad (7)$$

where $\lambda = \mathbb{E}(d_n^i)^2$, $\tau = \mathbb{E}(d_n^i)^4$, and $\kappa = (\tau - \lambda^2)$ for any $i = 1, \dots, N$

Zero-mean feedback term

Zero-mean term

Mean of the Hessian estimate

$$\begin{aligned}\mathbb{E} \left[\widehat{\mathbf{H}}_n \mid \mathcal{F}_n \right] &= \nabla^2 f(\theta_n) + \mathbb{E} \left[\Psi_n(\nabla^2 f(\theta_n)) \mid \mathcal{F}_n \right] + O(\delta_n^2) \\ &+ \mathbb{E} \left[\left(\frac{\xi_n^+ + \xi_n^- - 2\xi_n}{\delta_n^2} \right) \mid \mathcal{F}_n \right]\end{aligned}\tag{8}$$

Zero-mean

¹For any matrix P , $[P]_D$ refers to a matrix that retains only the diagonal entries of P and replaces all the remaining entries with zero

² $[P]_N$ to refer to a matrix that retains only the off-diagonal entries of P , while replaces all the diagonal entries with zero

Zero-mean feedback term

Zero-mean term

Mean of the Hessian estimate

$$\begin{aligned}\mathbb{E} \left[\widehat{\mathbf{H}}_n \mid \mathcal{F}_n \right] &= \nabla^2 f(\theta_n) + \mathbb{E} \left[\Psi_n(\nabla^2 f(\theta_n)) \mid \mathcal{F}_n \right] + O(\delta_n^2) \\ &+ \mathbb{E} \left[\left(\frac{\xi_n^+ + \xi_n^- - 2\xi_n}{\delta_n^2} \right) \mid \mathcal{F}_n \right]\end{aligned}\tag{8}$$

Zero-mean

Feedback term

$$\Psi_n(\mathbf{H}) = [\mathbf{M}_n]_D (d_n^T [\mathbf{H}]_N d_n) + [\mathbf{M}_n]_N (d_n^T [\mathbf{H}]_D d_n)\tag{9}$$

¹ For any matrix \mathbf{P} , $[\mathbf{P}]_D$ refers to a matrix that retains only the diagonal entries of \mathbf{P} and replaces all the remaining entries with zero

² $[\mathbf{P}]_N$ to refer to a matrix that retains only the off-diagonal entries of \mathbf{P} , while replaces all the diagonal entries with zero

Problem

Feedback term is function of current Hessian $\nabla^2 f$

Problem

Feedback term is function of current Hessian $\nabla^2 f$

Solution

Use \bar{H}_{n-1} as a proxy for $\nabla^2 f$

$$\hat{\Psi}_n = \Psi_n(\bar{H}_{n-1}) \quad (10)$$

Step-size optimization

Recall the Hessian recursion, $\bar{H}_n = (1 - b_n)\bar{H}_{n-1} + b_n(\hat{H}_n - \hat{\Psi}_n)$

Step-size optimization

Recall the Hessian recursion, $\bar{H}_n = (1 - b_n)\bar{H}_{n-1} + b_n(\hat{H}_n - \hat{\Psi}_n)$

Rewriting the Hessian recursion

$$\bar{H}_n = \sum_{i=0}^n \tilde{b}_i (\hat{H}_i - \hat{\Psi}_i) \quad (11)$$

Step-size optimization

Recall the Hessian recursion, $\bar{H}_n = (1 - b_n)\bar{H}_{n-1} + b_n(\hat{H}_n - \hat{\Psi}_n)$

Rewriting the Hessian recursion

$$\bar{H}_n = \sum_{i=0}^n \tilde{b}_i (\hat{H}_i - \hat{\Psi}_i) \quad (11)$$

Optimization problem for weights

$$\min_{\{\tilde{b}_i\}} \sum_{i=0}^n (\tilde{b}_i)^2 \delta_i^{-4}, \text{ subject to} \quad (12)$$

$$\tilde{b}_i \geq 0 \quad \forall i \text{ and } \sum_{i=0}^n \tilde{b}_i = 1 \quad (13)$$

Above optimization problem solution

$$\tilde{b}_i^* = \delta_i^4 / \sum_{j=0}^n \delta_j^4, i = 1, \dots, n \quad (14)$$

¹Step-size optimization is a relatively straightforward migration from Spall
2009

Above optimization problem solution

$$\tilde{b}_i^* = \delta_i^4 / \sum_{j=0}^n \delta_j^4, i = 1, \dots, n \quad (14)$$

Optimal weights for original Hessian recursion

$$b_i = \delta_i^4 / \sum_{j=0}^i \delta_j^4 \quad (15)$$

¹Step-size optimization is a relatively straightforward migration from Spall 2009

Improved Hessian Estimation for
Simultaneous Perturbation Stochastic
Approximation with 3 Simulations
(SPSA-3)

Function measurements

$$y_n^+ = f(\theta_n + \delta_n \Delta_n) + \xi_n^+, \quad y_n^- = f(\theta_n - \delta_n \Delta_n) + \xi_n^-$$

Gradient estimate

$$\hat{\nabla} f(\theta_n) = \begin{pmatrix} \frac{f(\theta_n + \delta_n \Delta_n) - f(\theta_n - \delta_n \Delta_n)}{2\delta_n \Delta_{n1}} + \frac{\xi_n^+ - \xi_n^-}{2\delta_n \Delta_{n1}} \\ \vdots \\ \frac{f(\theta_n + \delta_n \Delta_n) - f(\theta_n - \delta_n \Delta_n)}{2\delta_n \Delta_{nN}} + \frac{\xi_n^+ - \xi_n^-}{2\delta_n \Delta_{nN}} \end{pmatrix} \quad (16)$$

¹J. C. Spall (1992), "Multivariate stochastic approximation using a simultaneous perturbation gradient approximation," IEEE TAC.

Function measurements

$$y_n^{++} = f(\theta_n + \delta_n \Delta_n + \delta_n \widehat{\Delta}_n) + \xi_n^{++}, \quad y_n = f(\theta_n) + \xi_n,$$

$$y_n^{--} = f(\theta_n - \delta_n \Delta_n - \delta_n \widehat{\Delta}_n) + \xi_n^{--}$$

Hessian estimate \widehat{H}_n

$$\left(\widehat{H}_n\right)_{ij} = \left(\frac{y_n^{++} + y_n^{--} - 2y_n}{2\delta_n^2 \Delta_n^{(i)} \widehat{\Delta}_n^{(j)}}\right) \quad (17)$$

¹S. Bhatnagar, Prashanth L. A (2015) ,“Simultaneous perturbation Newton algorithms for simulation optimization,” JOTA.

Simplified Hessian estimate

$$\widehat{H}_n = \nabla^2 f(\theta_n) + \Psi_n(\nabla^2 f(\theta_n)) + O(\delta_n^2) + O(\delta_n^{-2}) \quad (18)$$

Feedback term

$$\Psi_n(H) = \frac{1}{2} P_n \left[\Delta_n^T H \Delta_n + \widehat{\Delta}_n^T H \widehat{\Delta}_n \right] + \widehat{N}_n^T H N_n + \widehat{N}_n^T H + H N_n \quad (19)$$

Where,

$$P_n = [1./\Delta_n][1./\Delta_n]^T, \quad N_n = \Delta_n[1./\Delta_n]^T - I_N, \quad \widehat{N}_n = \widehat{\Delta}_n[1./\widehat{\Delta}_n]^T - I_N$$

Lemma

(**Bias in Hessian estimate**) Under assumptions similar to those for 2SPSA and 2RDSA, we have a.s. that¹, for $i, j = 1, \dots, N$,

$$\left| \mathbb{E} \left[\widehat{H}_n(i, j) \middle| \mathcal{F}_n \right] - \nabla_{ij}^2 f(\theta_n) \right| = O(\delta_n^2) \quad (20)$$

Theorem

(**Strong Convergence of Hessian**) Under assumptions similar to those for 2SPSA and 2RDSA, we have that

$$\theta_n \rightarrow \theta^*, \bar{H}_n \rightarrow \nabla^2 f(\theta^*) \text{ a.s. as } n \rightarrow \infty$$

¹Here $\widehat{H}_n(i, j)$ and $\nabla_{ij}^2 f(\cdot)$ denote the (i, j) th entry in the Hessian estimate \widehat{H}_n and the true Hessian $\nabla^2 f(\cdot)$, respectively

Convergence analysis

Recall the Hessian recursion, $\bar{H}_n = (1 - b_n)\bar{H}_{n-1} + b_n(\hat{H}_n - \hat{\Psi}_n)$

Recall the Hessian recursion, $\bar{H}_n = (1 - b_n)\bar{H}_{n-1} + b_n(\hat{H}_n - \hat{\Psi}_n)$

Rewriting Hessian recursion as SA scheme

$$\begin{aligned}\bar{H}_n &= \bar{H}_{n-1} - b_n(\bar{H}_{n-1} - \hat{H}_n + \hat{\Psi}_n) \\ &= \bar{H}_{n-1} - b_n(\bar{H}_{n-1} - H^* + \hat{\Psi}_n - \Psi_n(H^*))\end{aligned}\quad (21)$$

Recall the Hessian recursion, $\bar{H}_n = (1 - b_n)\bar{H}_{n-1} + b_n(\hat{H}_n - \hat{\Psi}_n)$

Rewriting Hessian recursion as SA scheme

$$\begin{aligned}\bar{H}_n &= \bar{H}_{n-1} - b_n(\bar{H}_{n-1} - \hat{H}_n + \hat{\Psi}_n) \\ &= \bar{H}_{n-1} - b_n(\bar{H}_{n-1} - H^* + \hat{\Psi}_n - \Psi_n(H^*))\end{aligned}\quad (21)$$

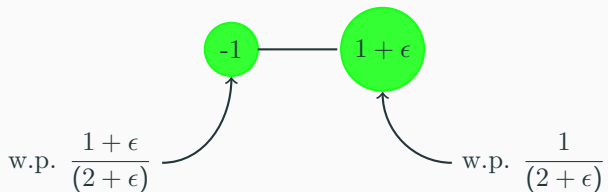
Theorem

(2SPSA-3-IH Quadratic case - Convergence rate) Let $b_n = b_0/n^r$, where $1/2 < r < 1$ and $0 < b_0 \leq 1$, $H^* = \nabla^2 f(\theta^*)$ and $\Lambda_k = \bar{H}_k - H^*$. Under noise-free setting, we have

$$\text{trace}[\mathbb{E}(\Lambda_n^T \Lambda_n)] = O(e^{-2b_0 n^{1-r}/(1-r)}) \quad (22)$$

Numerical Results

Asymmetric Bernoulli distribution¹



Uniform perturbations¹

$$d_n^i = \text{Unif}[-\eta, \eta], \eta > 0 \quad (23)$$

¹Prashanth L. A, Shalabh Bhatnagar, Michael Fu, Steve Marcus (2016), "Adaptive system optimization using random directions stochastic approximation," IEEE TAC.

Quadratic loss

$$f(\theta) = \theta^T A \theta + b^T \theta \quad (24)$$

Fourth-order loss

$$f(\theta) = \theta^T A^T A \theta + 0.1 \sum_{j=1}^N (A\theta)_j^3 + 0.01 \sum_{j=1}^N (A\theta)_j^4 \quad (25)$$

Additive Noise : $[\theta^T, 1]Z$, where $Z \approx \mathcal{N}(0, \sigma^2 I_{N+1 \times N+1})$

¹The implementation is available at <https://github.com/prashla/RDSA/archive/master.zip>

Normalized MSE (NMSE)

$$\|\theta_{\text{nend}} - \theta^*\|^2 / \|\theta_0 - \theta^*\|^2 \quad (26)$$

Normalized loss

$$f(\theta_{\text{nend}})/f(\theta_0) \quad (27)$$

Table 1: Normalized loss values for fourth-order objective (25) with noise: simulation budget = 10,000 and standard error from 500 replications shown after \pm

Noise parameter $\sigma = 0.1$		
	Regular	Improved Hessian estimation
2SPSA	0.132 ± 0.0267	0.104 ± 0.0355
2SPSA-3	0.0951 ± 0.0031	0.0594 ± 0.0014
2RDSA-Unif	0.115 ± 0.0214	0.0271 ± 0.0538
2RDSA-AsymBer	0.0471 ± 0.021	0.0099 ± 0.0014

¹ **Observation 1:** Schemes with improved Hessian estimation performs better than their respective regular schemes

² **Observation 2:** 2RDSA-IH-AsymBer is performing the best overall

Table 2: NMSE values for quadratic objective (24) with noise: simulation budget = 10,000 and standard error from 500 replications shown after \pm

Noise parameter $\sigma = 0.1$		
	Regular	Improved Hessian estimation
2SPSA	0.9491 ± 0.0131	0.5495 ± 0.0217
2SPSA-3	0.8378 ± 0.0179	0.1045 ± 0.0005
2RDSA-Unif	1.0073 ± 0.0140	0.1953 ± 0.0095
2RDSA-AsymBer	0.1667 ± 0.0095	0.0324 ± 0.0007

¹ **Observation 1:** Schemes with improved Hessian estimation performs better than their respective regular schemes

² **Observation 2:** 2RDSA-IH-AsymBer is performing the best overall

- Proposed generalised RDSA algorithm + Improved Hessian estimation scheme
- Improved Hessian estimation scheme for 2SPSA-3 algorithm
- 2RDSA-IH, 2SPSA-3 requires only 75% of the simulation cost per-iteration for 2SPSA, 2SPSA-IH

- To improve rate of convergence of first-order methods by incorporating ideas of momentum descent, Hessian-free methods, and conjugate methods
- To develop stochastic approximation versions of quasi-Newton schemes given by the Broyden family
- To explore applications of these methods to the design of reinforcement learning algorithms
- Deterministic perturbations for Second-order methods

Thank You