

SARAH: Stochastic Recursive grAdient algorithM

Pang-Chun Chung, Hung-Ruei Wu

Department of Computer Science, National Yang Ming Chiao Tung University, Hsinchu, Taiwan



Motivation

- **Limitation of SVRG:** In the inner loop, SVRG only uses the true gradient as a historical reference. This limited information leads to *unstable updates* within the inner loop.
- **Limitation of SAG/SAGA:** SAG and SAGA utilize all historical gradient information, resulting in *high storage requirements*, which can be impractical for large-scale problems.
- **Proposed Solution - SARAH:** To address these limitations, we propose a *recursive update scheme* in SARAH. Our method combines the advantages of SVRG and SAG/SAGA by:
 - ▷ Ensuring **stability** through recursive updates within the inner loop.
 - ▷ Avoiding excessive **storage costs** while maintaining efficiency.

Finite-Sum Optimization Problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} P(\mathbf{w}) := \frac{1}{n} \sum_{i \in [n]} f_i(\mathbf{w}),$$

where f_i is convex with a Lipschitz continuous gradient, $i \in [n] := \{1, \dots, n\}$, and we assume that the optimal solution \mathbf{w}^* exists.

Remark: This optimization problem forms the foundation for many algorithms, including SVRG and our proposed SARAH.

SARAH: Stochastic Recursive Gradient Algorithm

Input: Learning rate $\eta > 0$, inner loop size m , initial point \mathbf{w}_0
Output: Optimized solution \mathbf{w}_s

Algorithm 1: SARAH

```
Initialize  $\tilde{\mathbf{w}}_0$ 
For  $s = 1, 2, \dots$ 
     $\mathbf{w}_0 \leftarrow \tilde{\mathbf{w}}_{s-1}$ 
     $\mathbf{v}_0 \leftarrow \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}_0)$ 
     $\mathbf{w}_1 \leftarrow \mathbf{w}_0 - \eta \mathbf{v}_0$ 
    For  $t = 1, \dots, m-1$ 
        Sample  $i_t$  uniformly at random from  $[n]$ 
         $\mathbf{v}_t \leftarrow \nabla f_{i_t}(\mathbf{w}_t) - \nabla f_{i_t}(\mathbf{w}_{t-1}) + \mathbf{v}_{t-1}$ 
         $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \mathbf{v}_t$ 
    End For
    Set  $\tilde{\mathbf{w}}_s \leftarrow \mathbf{w}_t$  with  $t$  chosen uniformly at random from  $\{0, 1, \dots, m\}$ 
End For
```

Assumptions

- **Assumption 1 (L-smooth).** Each $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, $i \in [n]$, is L -smooth, i.e.,
$$\|\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}')\| \leq L \|\mathbf{w} - \mathbf{w}'\|, \quad \forall \mathbf{w}, \mathbf{w}' \in \mathbb{R}^d.$$
- **Assumption 2a (μ -strongly convex).** The function $P : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex, i.e.,
$$P(\mathbf{w}) \geq P(\mathbf{w}') + \nabla P(\mathbf{w}')^T (\mathbf{w} - \mathbf{w}') + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}'\|^2.$$
- **Assumption 2b.** Each $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, $i \in [n]$, is strongly convex with $\mu > 0$. Strong convexity implies:
$$2\mu[P(\mathbf{w}) - P(\mathbf{w}^*)] \leq \|\nabla P(\mathbf{w})\|^2, \quad \forall \mathbf{w} \in \mathbb{R}^d.$$
- **Assumption 3 (Convexity).** Each $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, $i \in [n]$, is convex, i.e.,
$$f_i(\mathbf{w}) \geq f_i(\mathbf{w}') + \nabla f_i(\mathbf{w}')^T (\mathbf{w} - \mathbf{w}'), \quad \forall i \in [n].$$

Performance Metrics

- **ϵ -accurate Solution:** The solution \mathbf{w}_T satisfies $\|\nabla P(\mathbf{w}_T)\|^2 \leq \epsilon$, ensuring the gradient norm meets the desired accuracy level.
- **Full Gradient Calculation:** The analysis is based on the total number of stochastic gradient evaluations, bounding the iterations required to achieve the specified accuracy.

Comparison of SARAH and SVRG on Convex Problems

Method	Problem Type	Complexity
SARAH	Strongly Convex	$O((n + \kappa) \log(1/\epsilon))$
SVRG	Strongly Convex	$O((n + \kappa) \log(1/\epsilon))$
SARAH	General Convex	$O((n + 1/\epsilon) \log(1/\epsilon))$
SVRG	General Convex	$O(n + (n/\epsilon))$

SARAH+: A Practical Variant of SARAH

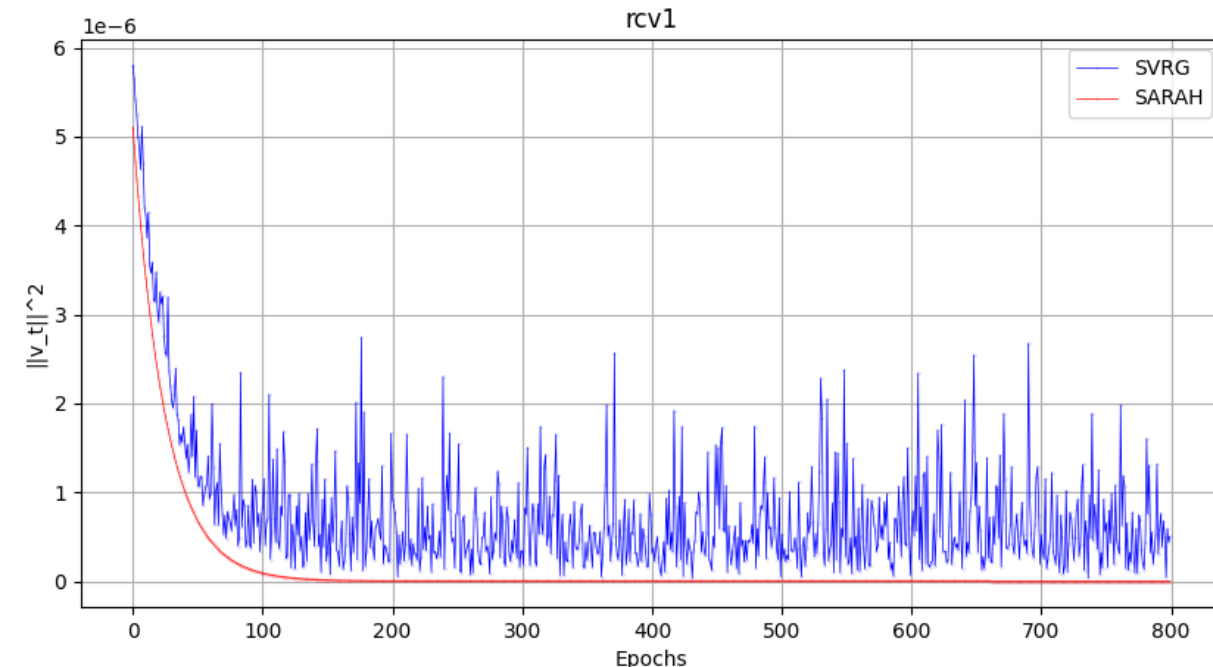
Input: Learning rate $\eta > 0$, $0 < \gamma \leq 1$, and the maximum inner loop size m
Output: Optimized solution $\tilde{\mathbf{w}}_s$

Algorithm 2: SARAH+

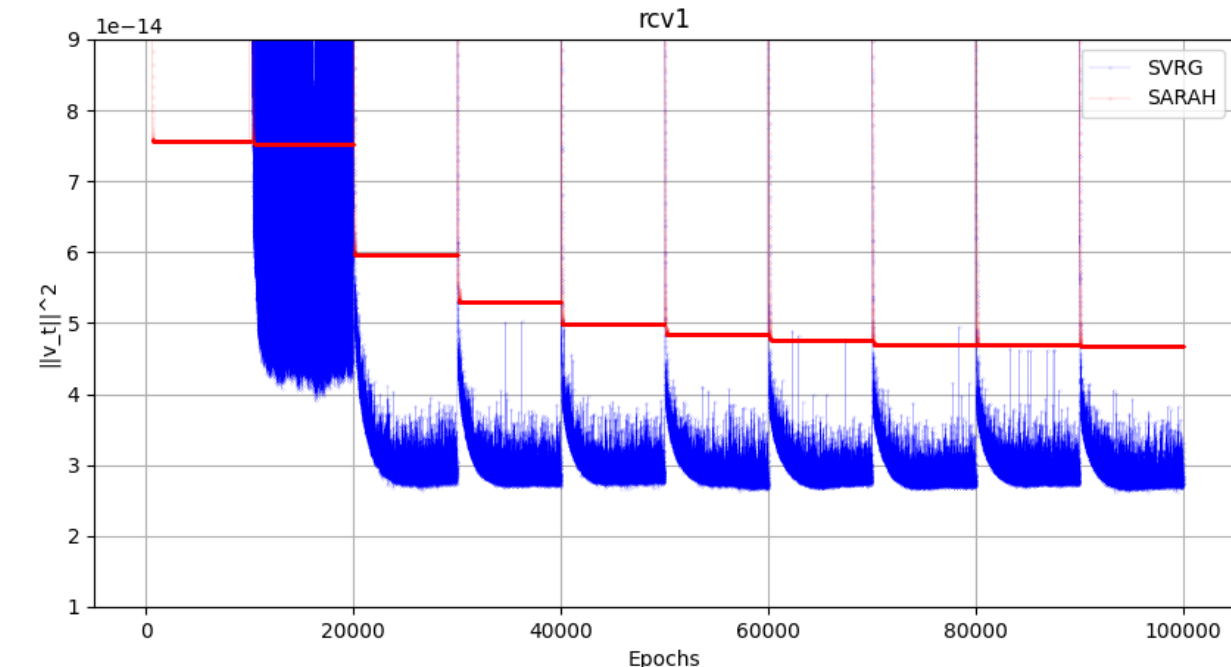
```
Initialize  $\tilde{\mathbf{w}}_0$ 
For  $s = 1, 2, \dots$ 
     $\mathbf{w}_0 \leftarrow \tilde{\mathbf{w}}_{s-1}$ 
     $\mathbf{v}_0 \leftarrow \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}_0)$ 
     $\mathbf{w}_1 \leftarrow \mathbf{w}_0 - \eta \mathbf{v}_0$ 
     $t \leftarrow 1$ 
    While  $\|\mathbf{v}_{t-1}\|^2 > \gamma \|\mathbf{v}_0\|^2$  and  $t < m$ 
        Sample  $i_t$  uniformly at random from  $[n]$ 
         $\mathbf{v}_t \leftarrow \nabla f_{i_t}(\mathbf{w}_t) - \nabla f_{i_t}(\mathbf{w}_{t-1}) + \mathbf{v}_{t-1}$ 
         $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \mathbf{v}_t$ 
         $t \leftarrow t + 1$ 
    End While
    Set  $\tilde{\mathbf{w}}_s \leftarrow \mathbf{w}_t$ 
End For
```

Numerical Result I

The following figures show the stochastic gradient of SVRG and SARAH.



(a) Single Inner Loop



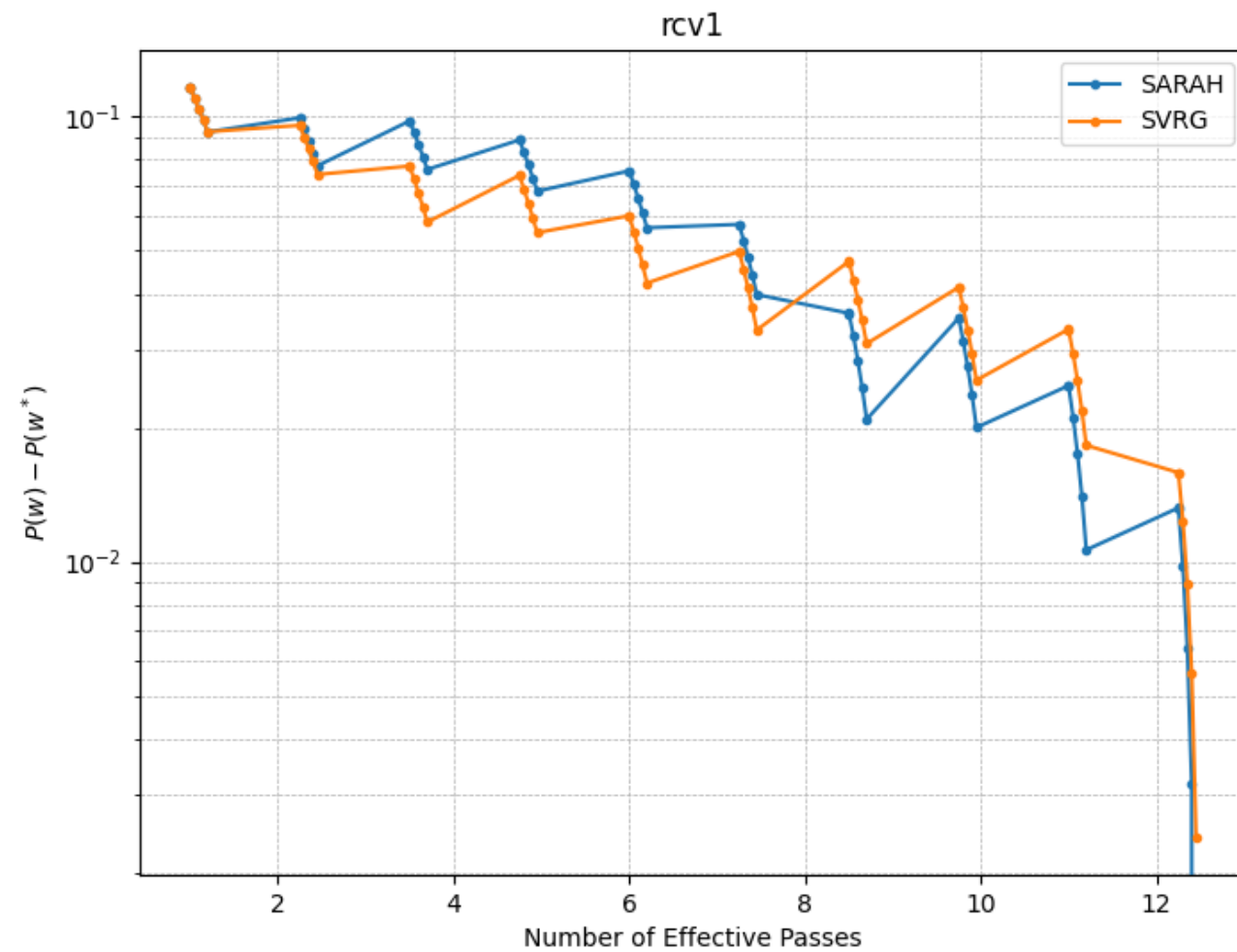
(b) Multiple Inner Loop

Key Observations:

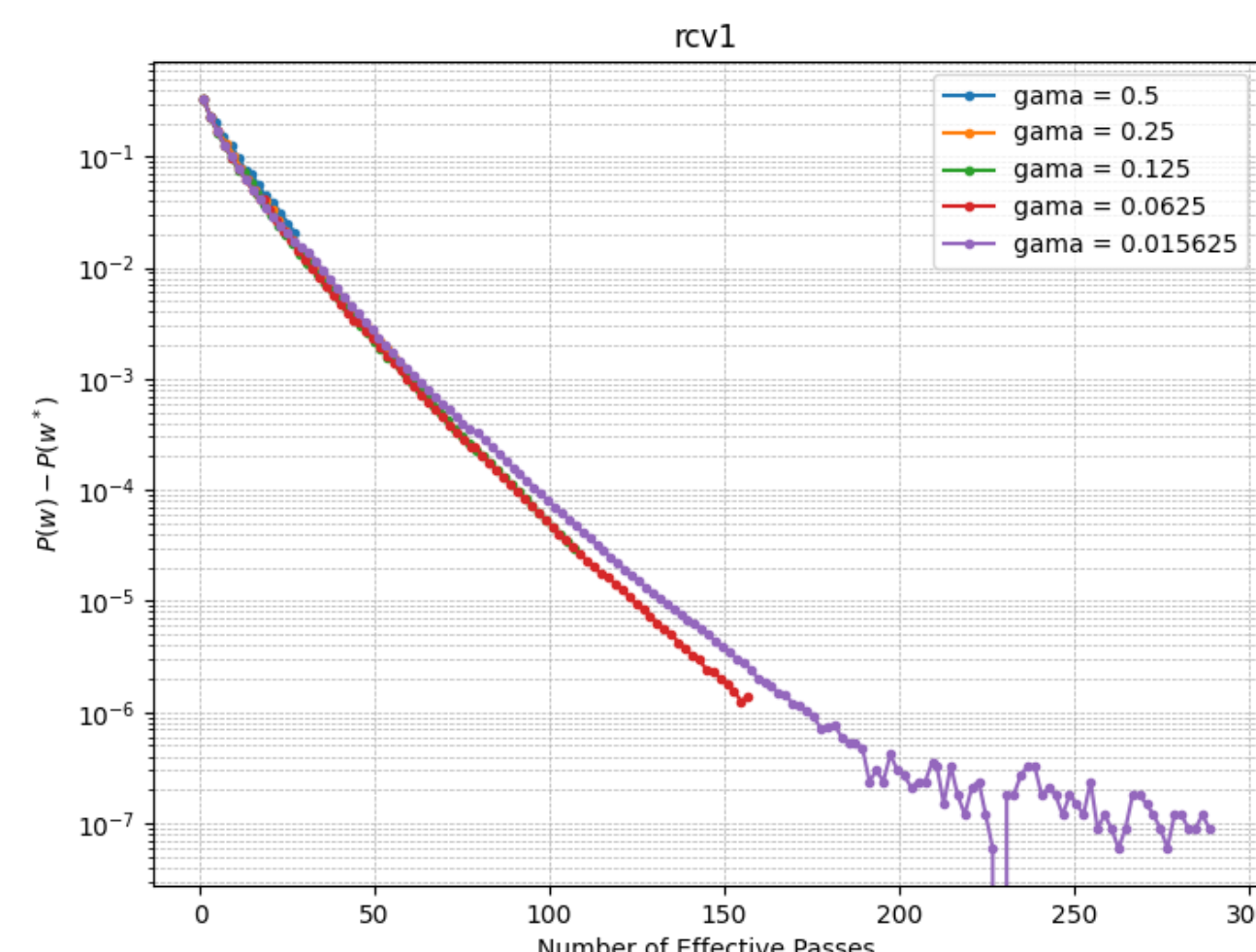
- SVRG shows oscillations due to large fluctuations in stochastic gradients.
- SARAH is more stable and converges smoothly.

Numerical Result II

The figure below compares the performance of SVRG and SARAH in terms of loss residuals.



(a) Comparison of SVRG and SARAH



(b) Comparison SARAH+ with different γ

Key Observations:

- **Left Figure (SVRG and SARAH):** As mentioned in the paper, SARAH initially does not perform as well as SVRG. However, with sufficient effective passes, SARAH slightly outperforms SVRG in terms of convergence stability and final accuracy.
- **Right Figure (SARAH+ with gamma):** The performance of SARAH+ depends on the choice of gamma values. The most appropriate range for γ appears to be approximately between $\frac{1}{8}$ and $\frac{1}{16}$.

Conclusion

- A new **variance reducing** algorithm has been proposed.
- The algorithm achieves a **linear convergence** rate for the inner loop.
- The convergence rate's **constant is smaller**, improving efficiency.
- SARAH+ provides an explicit **stopping criterion** for practical implementations.

References

- Johnson, Rie, and Tong Zhang, "Accelerating Stochastic Gradient Descent Using Predictive Variance Reduction," NeurIPS, 2013.
- Nguyen, Lam M., Jie Liu, Katya Scheinberg, and Martin Takáč, "SARAH: A Novel Method for Machine Learning Problems Using Stochastic Recursive Gradient," ICML, 2017.