

# Riemannian stochastic variance reduced gradient on Grassmann manifold

Hiroyuki Kasai\*, Hiroyuki Sato<sup>§</sup> and Bamdev Mishra<sup>†</sup>

\*The University of Electro-Communications, Japan <sup>§</sup>Tokyo university of Science, Japan

<sup>†</sup>Amazon Development Centre India, India

## Summary (our contributions)

- Address stochastic gradient descent (SGD) algorithm for empirical risk minimization problem as

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n f_i(w).$$

- Paritularly, **structured** problems on **manifolds**, i.e.,  $w \in \mathcal{M}$  [1].
- Propose **Riemannian SVRG (R-SVRG)**.
- Extend **SVRG** in the **Euclidean** into **Riemannian manifolds**.
- **Global and Local convergence** analyses for (i) **exponential mapping and parallel translation**, and (ii) **retraction and vector transport**.
- Show effectiveness of R-SVRG from numerical comparisons.

## Stochastic variance reduced gradient (SVRG) [2]

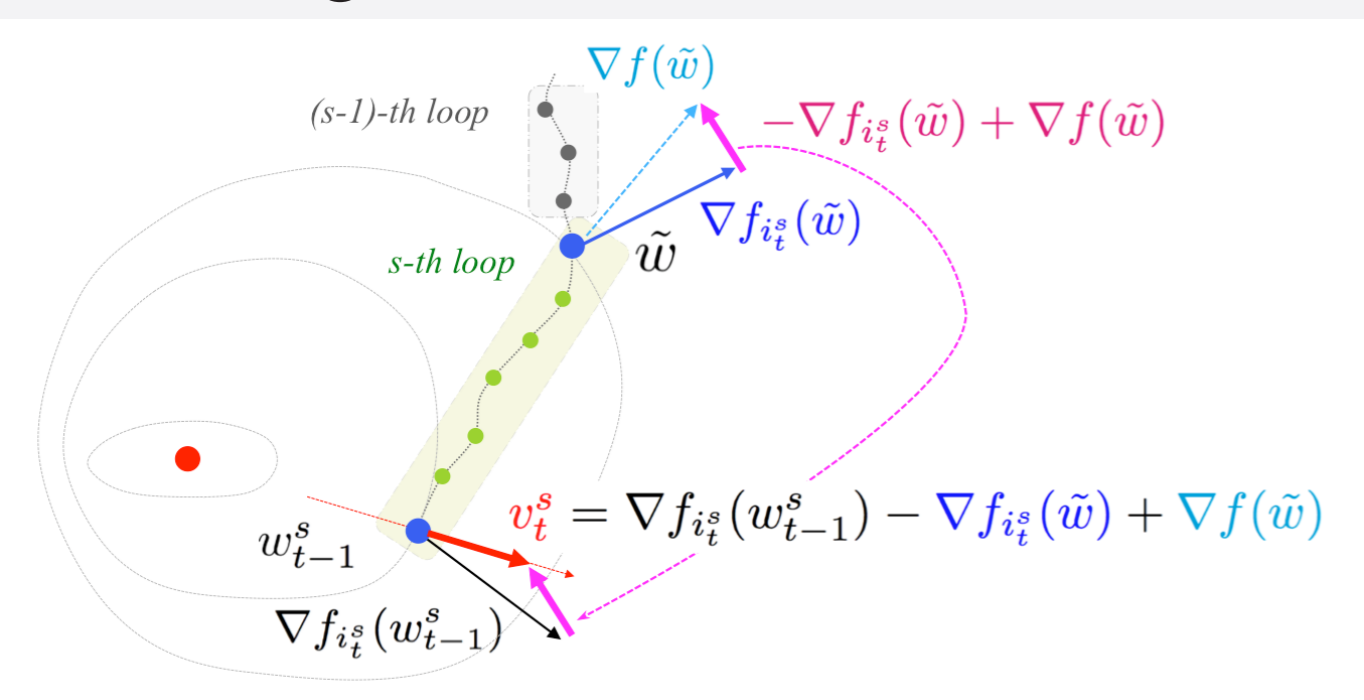
- Stochastic gradient descent (SGD)
  - High scalability to large-scale data because per-iteration complexity does not depend on  $n$ .
  - But, slow convergence property (sub-linear) with decaying stepsize
- Stochastic variance reduced gradient (SVRG)

### Basic idea

- Reduce the **variance** of stochastic gradients.
- **No need** to store all gradients.
- **Linear** convergence rate

### Basic strategy

- **Periodically**, calculate and store a **full gradient**.
- Every iteration **adjusts** a stochastic gradient by the latest full gradient to.



### SVRG pseudo-code

- 1: Initial iterate  $w_0^s \in \mathcal{M}$ .
- 2: **for**  $s = 1, 2, \dots$  (outer loop) **do**
- 3: Store  $\tilde{w} = w_{t-1}^{s-1}$ .
- 4: Store  $\nabla f(\tilde{w})$ .
- 5: **for**  $t = 1, \dots, m_s$  (inner) **do**
- 6: Calculate modified stochastic gradient  $v_t^s = \nabla f_{i_t^s}(w_{t-1}^s) - \nabla f_{i_t^s}(\tilde{w}) + \nabla f(\tilde{w})$ .
- 7: Update  $w_t^s = w_{t-1}^s - \alpha v_t^s$ .
- 8: **end for**
- 9: set  $\tilde{w} = w_t^s$  for randomly chosen  $t$ .
- 10: **end for**

## Riemannian stochastic variance reduced gradient (R-SVRG) [3]

- Address various structured problems
  - E.g., PCA, matrix completion, subspace tracking, spectral clustering.
- Focus on the **Grassmann manifold  $\text{Gr}(r, d)$** .
  - The set of  $r$ -dimensional linear subspaces in  $\mathbb{R}^d$ .
  - Can be **generalized** to other **compact Riemannian manifolds**.

### Notations

	SVRG	R-SVRG
Model parameter	$w_{t-1}^s \in \mathbb{R}^n$	$U_{t-1}^s \in \text{Gr}(r, d)$
Edge point of outer loop	$\tilde{w} \in \mathbb{R}^n$	$\tilde{U} \in \text{Gr}(r, d)$
Stochastic gradient	$\nabla f_{i_t^s}(w_{t-1}^s) \in \mathbb{R}^n$	$\text{grad} f_{i_t^s}(U_{t-1}^s) \in T_{U_{t-1}^s} \text{Gr}(r, d)$
Modified stochastic gradient	$v_t^s \in \mathbb{R}^n$	$\xi_t^s \in T_{U_{t-1}^s} \text{Gr}(r, d)$

### Algorithm basis

- Build on Riemannian stochastic gradient [4].
- Straightforward modification of stochastic gradient
  - Extend SVRG case:  $v_t^s = \nabla f_{i_t^s}(w_{t-1}^s) - \nabla f_{i_t^s}(\tilde{w}) + \nabla f(\tilde{w})$

$$\xi_t^s = \text{grad} f_{i_t^s}(U_{t-1}^s) - \text{grad} f_{i_t^s}(\tilde{U}) + \text{grad} f(\tilde{U})$$

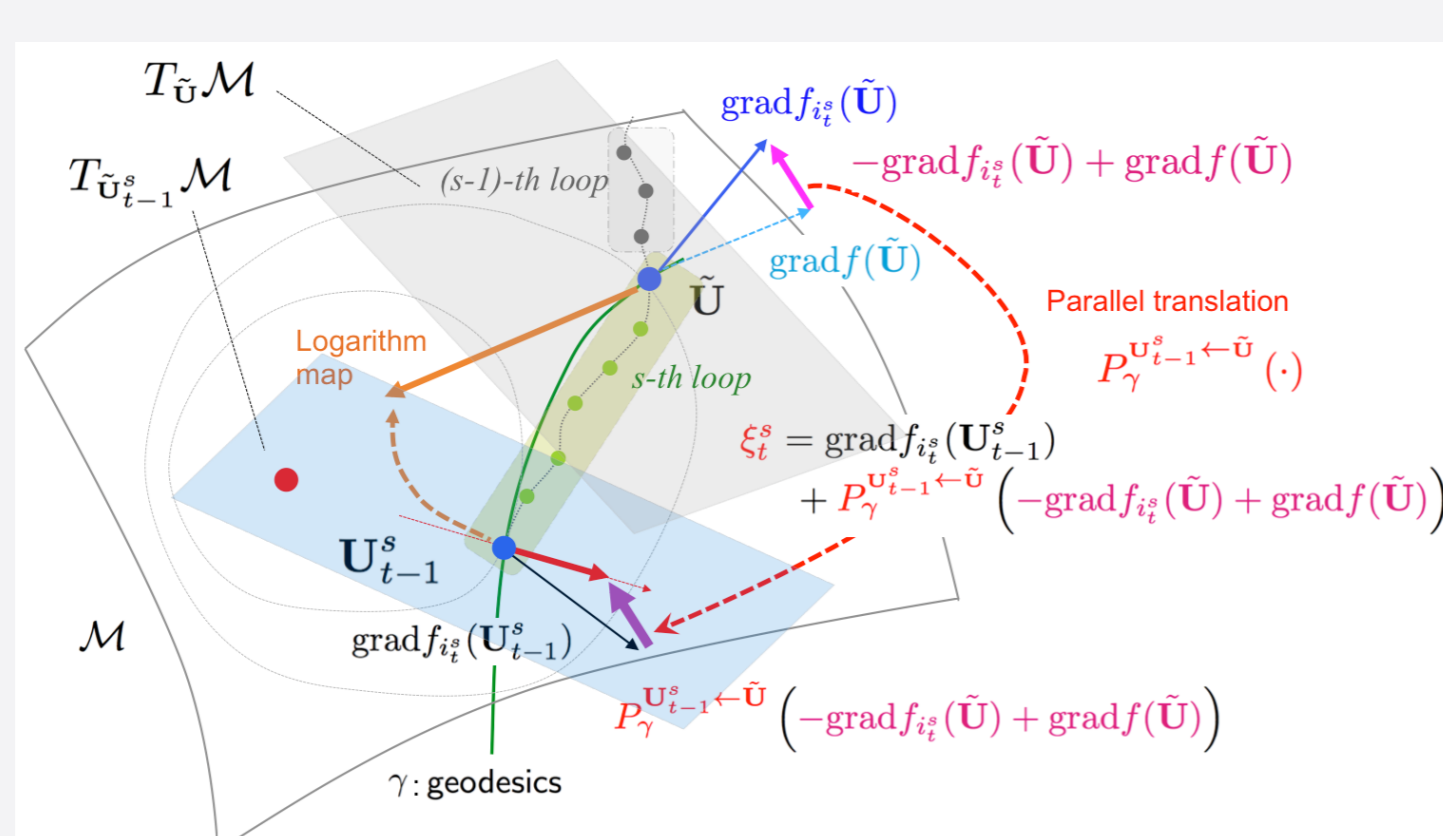
- **Meaningless** because manifolds are **not vector space**.

### Proposed modification

- **Transport** vectors at  $\tilde{U}$  into the current tangent space at  $U_{t-1}^s$  by **parallel translation** or **vector transport**, then add them.

$$\xi_t^s = \text{grad} f_{i_t^s}(U_{t-1}^s) + \mathcal{T} \left( -\text{grad} f_{i_t^s}(\tilde{U}) + \text{grad} f(\tilde{U}) \right)$$

parallel-translation / vector transport operator



## Convergence analyses

- Global convergence analysis with **decaying step-sizes**.
  - Guarantee that the iteration **globally converges** to a **critical point** starting from any initialization point.
- Local convergence rate analysis under **fixed step-size**.
  - Consider the rate in **neighborhood of a local minimum**.
  - Assume that **Lipschitz smoothness** and **lower bound of Hessian** hold only in this neighborhood.
  - Exponential mapping and parallel translation

$$\mathbb{E}[(\text{dist}(\tilde{U}^s, U^*))^2] \leq \frac{4(1 + 8m\alpha^2\beta^2)}{\alpha m(\sigma - 14\eta\beta^2)} \mathbb{E}[(\text{dist}(\tilde{U}^{s-1}, U^*))^2].$$

- Retraction and vector transport

$$\mathbb{E}[(\text{dist}(\tilde{U}^s, U^*))^2] \leq \frac{4(1 + 16m\alpha^2(\beta^2 + C^2a^2))}{\alpha m(\sigma - 28\alpha(\beta^2 + C^2a^2))} \mathbb{E}[(\text{dist}(\tilde{U}^{s-1}, U^*))^2],$$

where  $a$  and  $C$  are some constants satisfying  $\|\mathcal{T}_\eta(\xi) - P_\eta(\xi)\| \leq a\|\xi\|\|\eta\|$  and  $\|\text{grad} f_i(U)\| \leq C$  for  $U$  sufficiently close to the local minimum and for  $\xi, \eta \in T_U \text{Gr}(r, d)$ .

## Proof sketch in exponential mapping and parallel translation case

- Obtain below by assuming the **smallest eigenvalue  $\sigma$**  of Hessian of  $f$  as
 
$$f(z) \geq f(w) + \langle \text{Exp}_w^{-1}(z), \text{grad} f(w) \rangle_w + \frac{\sigma}{2} \|\text{Exp}_w^{-1}(z)\|_w^2, \quad w, z \in \mathcal{U}. \quad (1)$$

- Obtain the **variance of  $\xi_t^s$**  from  $\beta$ -Lipschitz continuity as

$$\mathbb{E}_{i_t^s}[\|\xi_t^s\|^2] \leq \beta^2(14(\text{dist}(w_{t-1}^s, w^*))^2 + 8\text{dist}(\tilde{w}^{s-1}, w^*))^2) \quad (2)$$

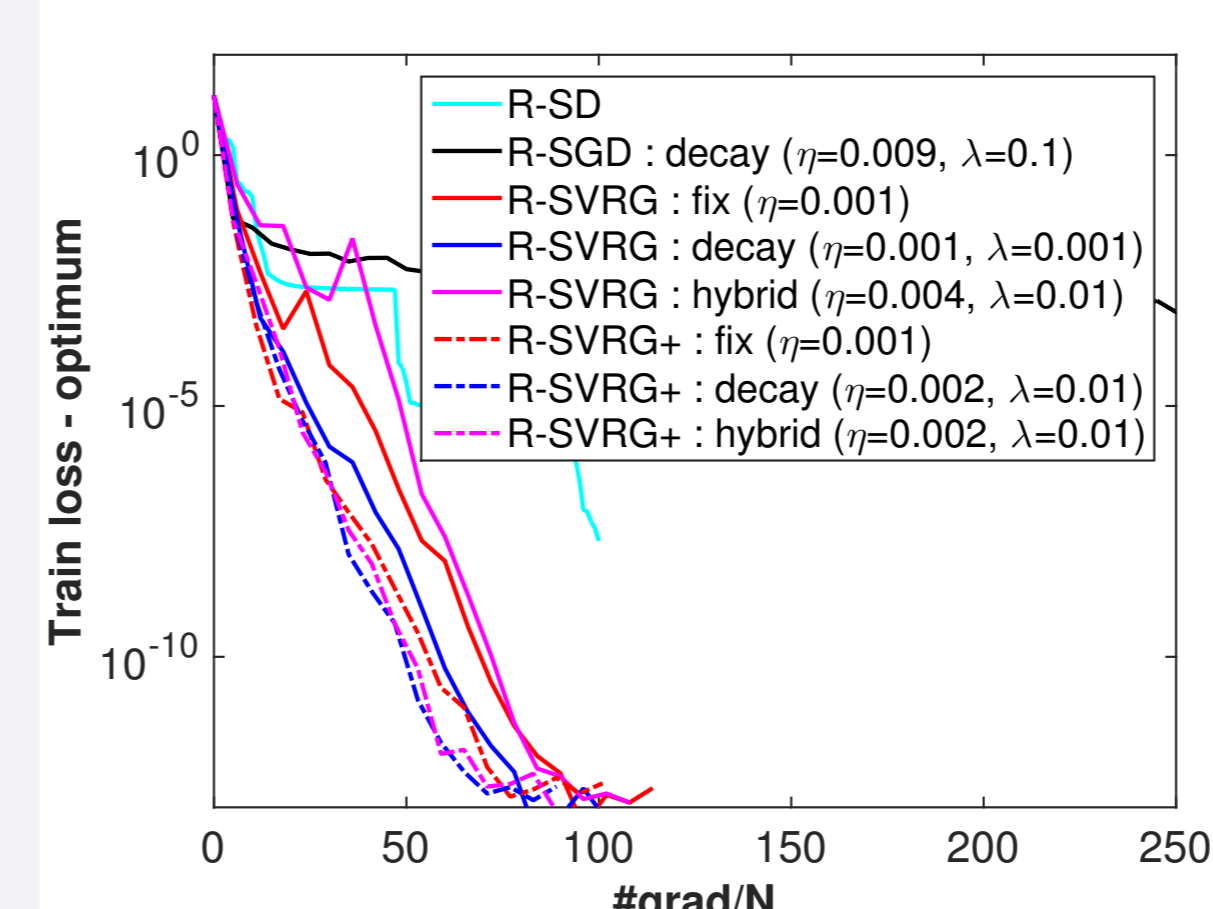
- Obtain the expectation of the **decrease of the distance** to the solution in the **inner iteration** from the lemma for a **geodesic triangle** in an Alexandrov space as

$$\mathbb{E}_{i_t^s} [(\text{dist}(U_t^s, U^*))^2 - (\text{dist}(U_{t-1}^s, U^*))^2] \leq \mathbb{E}_{i_t^s} [(\text{dist}(U_{t-1}^s, U^*))^2 + 2\eta \langle \text{grad} f(U_{t-1}^s), \text{Exp}_{U_{t-1}^s}^{-1}(U^*) \rangle_{U_{t-1}^s}]. \quad (3)$$

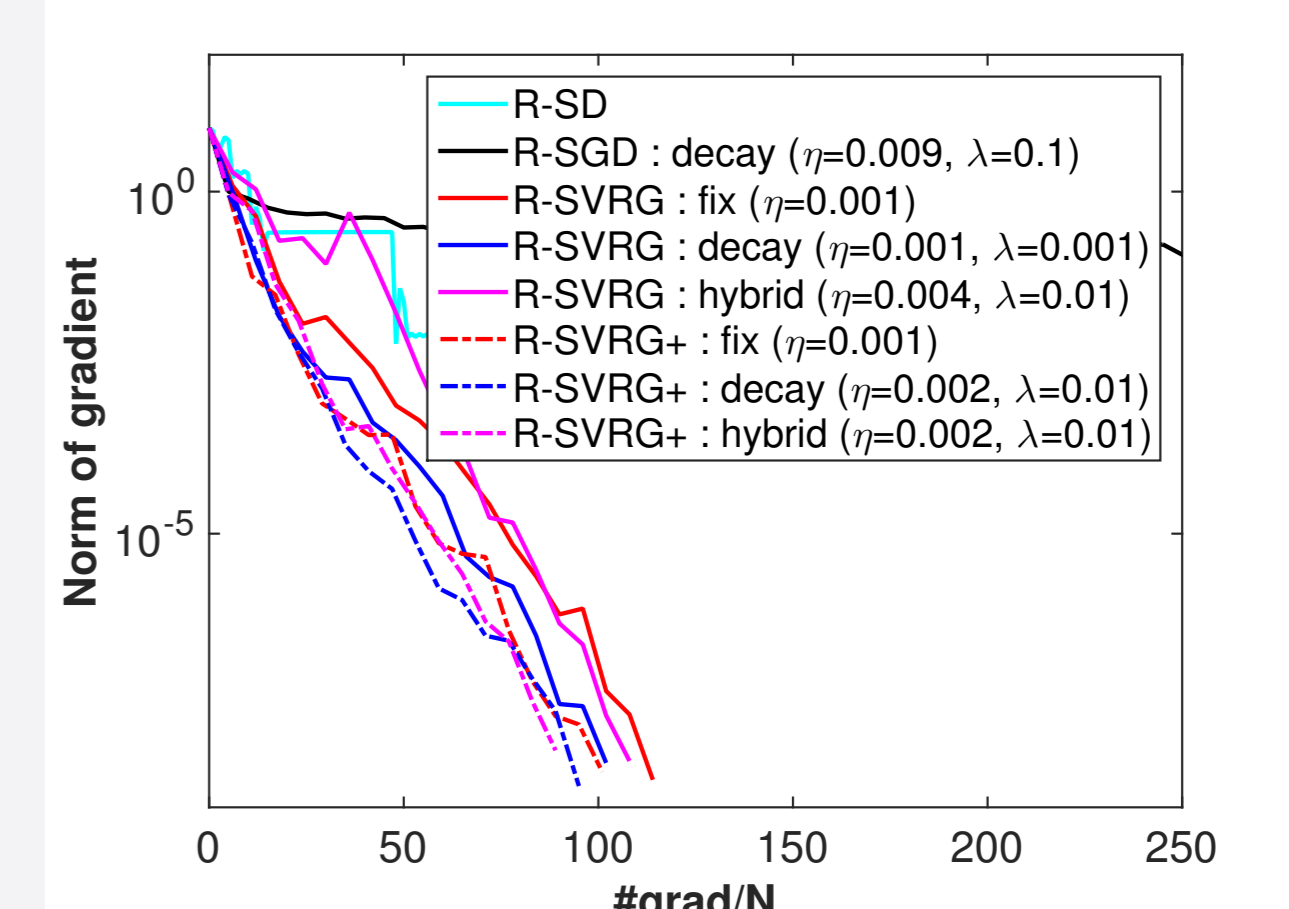
- Putting (1)&(2) into (3) with summing over the inner loop finally yields the **decrease of the distance** to the solution in the **outer iteration**.

## Numerical comparisons

- Evaluate the case of exponential mapping and parallel translation.
  - The tools in the Grassmann manifold have closed-form expressions [1].
- PCA problem (Synthetic dataset:  $n = 10000$ ,  $d = 20$ , and  $r = 5$ )

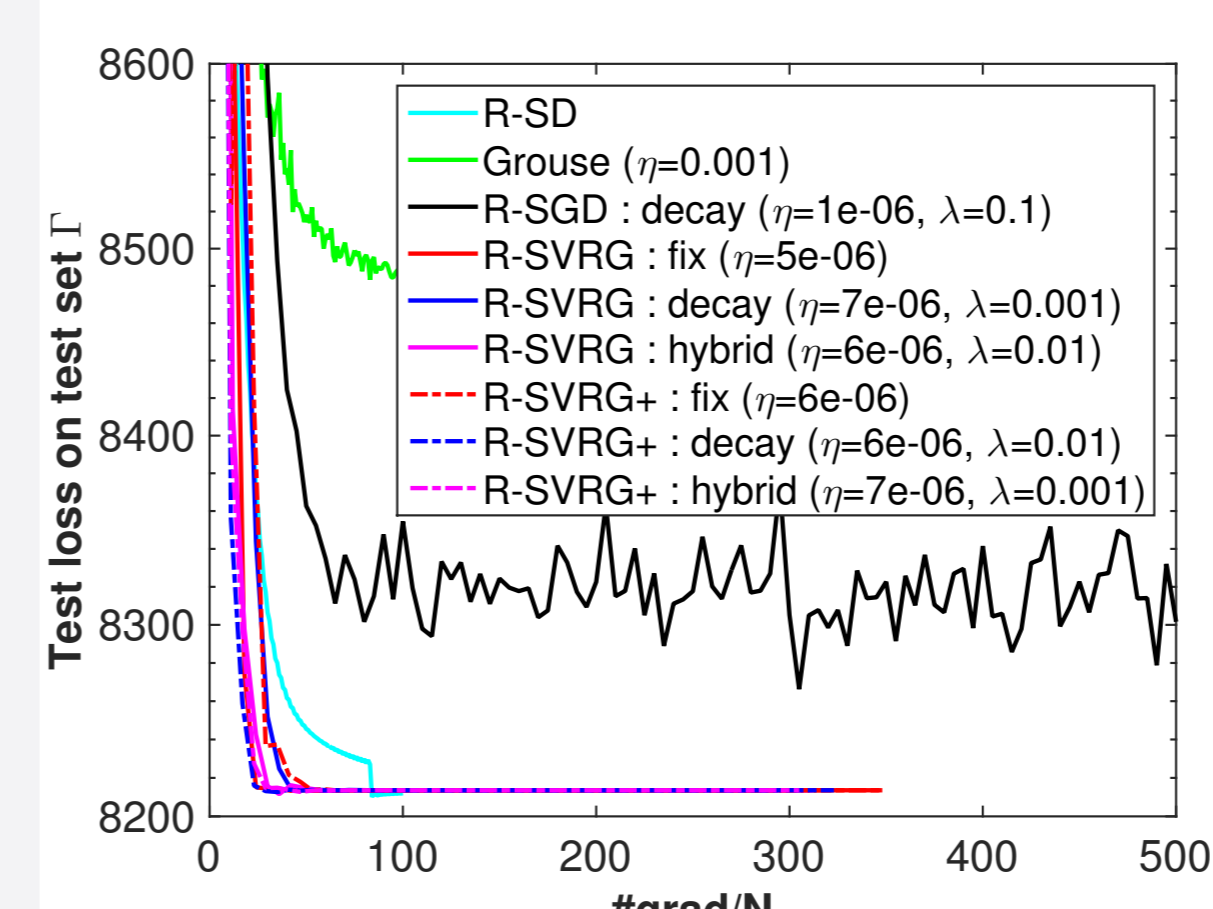


(a) Optimality gap.

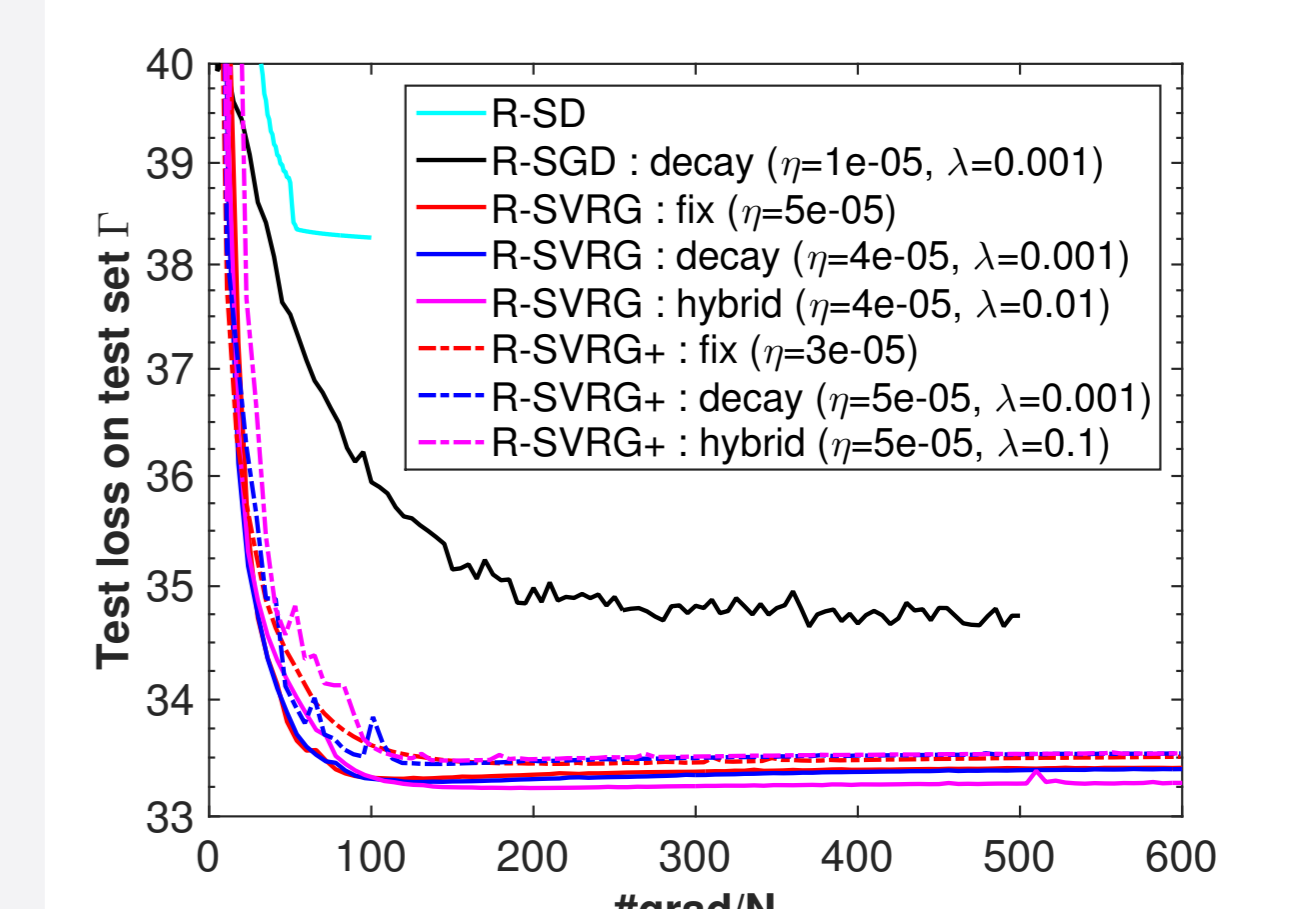


(b) Norm of gradient.

- Matrix completion problem (Jester and MovieLens10M)



(a) Jester dataset.



(b) MovieLens10M dataset.

## Reference

- [1] P.-A. Absil, R. Mahony, and R. Sepulchre. Optimization Algorithms on Matrix Manifolds. Princeton University Press, Princeton, NJ, 2008.
- [2] R. Johnson and T. Zhang. "Accelerating stochastic gradient descent using predictive variance reduction," In NIPS, pages 315-323, 2013.
- [3] H. Kasai, H. Sato, and B. Mishra, "Riemannian stochastic variance reduced gradient on Grassmann manifold," arXiv preprint: arXiv:1605.07367, 2016.
- [4] S. Bonnabel. Stochastic gradient descent on Riemannian manifolds. IEEE Trans. on Automatic Control, 58(9):2217-2229, 2013.