

RSN: Randomized Subspace Newton

Robert M. Gower¹ Dmitry Kovalev² Felix Lieder³ Peter Richtárik²

¹LTCI, Télécom Paris, Institut Polytechnique de Paris ²KAUST

³Heinrich-Heine-Universität Düsseldorf University



1. High Dimensional Optimization

Consider the optimization problem

$$x_* = \arg \min_{x \in \mathbb{R}^d} f(x), \quad (1)$$

where $f: \mathbb{R}^d \mapsto \mathbb{R}$ is C^2 and d is very big. This arises in training ML models with a very large number of parameters, or when data is high dimensional and acquiring data is expensive/hard.

Example: genomics, seismology, neurology and high resolution sensors in medicine.

Notation:

- Gradient & Hessian: $g(x) := \nabla f(x)$ & $\mathbf{H}(x) := \nabla^2 f(x)$
- Level set: $\mathcal{Q} := \{x \in \mathbb{R}^d : f(x) \leq f(x_0)\}$
- Hessian inner product: $\langle u, v \rangle_{\mathbf{H}(x)} := \langle \mathbf{H}(x)u, v \rangle$

2. Assumptions (New)

Assumption 1: Gradient invariance:

$$g(x) \in \text{Range}(\mathbf{H}(x)) \quad \text{for all } x \in \mathbb{R}^d. \quad (2)$$

Assumption 2: f is \hat{L} -smooth and $\hat{\mu}$ -convex relative to its Hessian. That is, there exist $\hat{L} \geq \hat{\mu} > 0$ such that for all $x, y \in \mathcal{Q}$:

$$f(x) \leq \underbrace{f(y) + \langle g(y), x - y \rangle + \frac{\hat{L}}{2} \|x - y\|_{\mathbf{H}(y)}^2}_{:=T(x,y)}, \quad (3)$$

$$f(x) \geq f(y) + \langle g(y), x - y \rangle + \frac{\hat{\mu}}{2} \|x - y\|_{\mathbf{H}(y)}^2. \quad (4)$$

This is a weak assumption since:

$$L\text{-smoothness} \Rightarrow c\text{-stability [1]} \Rightarrow \hat{L}\text{-smoothness}$$

$$\mu\text{-convexity} \Rightarrow \hat{\mu}\text{-convexity}$$

Example: Both assumptions hold for smooth generalized linear models with L_2 regularization.

3. Newton's Method

Newton's method applied to problem (1) has the form

$$x_{k+1} = x_k - \gamma \cdot \mathbf{H}^\dagger(x_k)g(x_k),$$

where

- $\gamma > 0$ is the stepsize
- $\mathbf{H}^\dagger(x_k)$ is the Moore-Penrose pseudoinverse of $\mathbf{H}(x_k)$

Pros: Can handle curvature, invariant to coordinate transformations

Cons: Cost of each iteration is very high: $\mathcal{O}(d^3)$

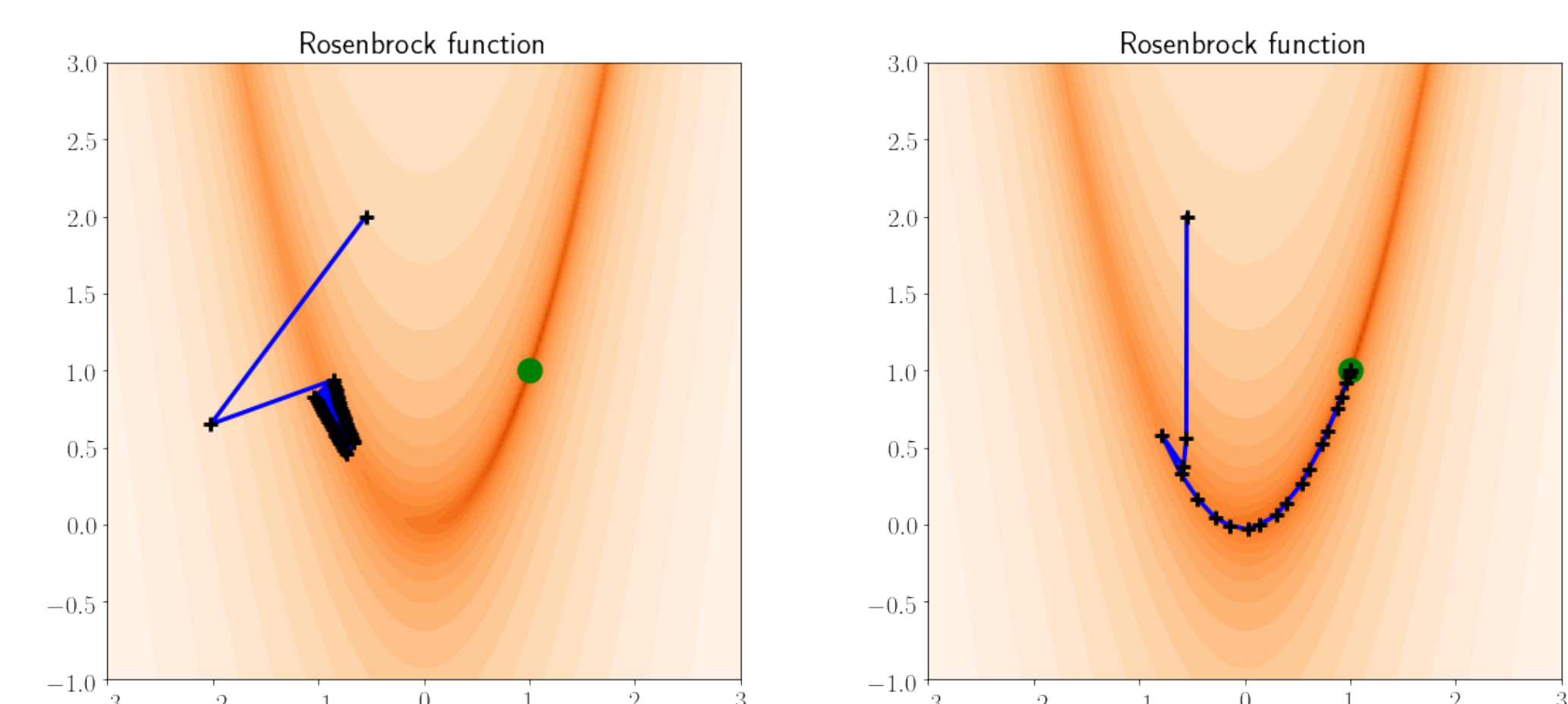


Figure: Gradient descent (left) and Newton's method (right) 50 iterations.

4. Sketching and Dimension Reduction

Let $\mathbf{S} \in \mathbb{R}^{d \times s}$ be a random matrix drawn from $\mathbf{S} \sim \mathcal{D}$.

$$\left. \begin{array}{c} \mathbf{S}^\top \\ \mathbf{x} \end{array} \right\} \in \mathbb{R}^d = \left. \begin{array}{c} \mathbf{S}^\top \mathbf{x} \\ \end{array} \right\} \in \mathbb{R}^s$$

Assumption 3: With probability 1, the sketching matrix \mathbf{S} satisfies:

$$\text{Null}(\mathbf{S}^\top \mathbf{H}(x) \mathbf{S}) = \text{Null}(\mathbf{S}), \quad \forall x \in \mathcal{Q}. \quad (5)$$

5. Randomized Subspace Newton

Algorithm 1 RSN: Randomized Subspace Newton

- 1: **input:** $x_0 \in \mathbb{R}^d$
- 2: **parameters:** \mathcal{D} = distribution over random matrices
- 3: **for** $k = 0, 1, 2, \dots$ **do**
- 4: Sample a fresh sketching matrix: $\mathbf{S}_k \sim \mathcal{D}$
- 5: $x_{k+1} = x_k - \frac{1}{\hat{L}} \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top g(x_k)$
- 6: **end for**
- 7: **output:** last iterate x_k

Computation of sketched Newton direction:

$$g(x) \xrightarrow{\mathbf{S}^\top} \mathbf{S}^\top g(x) \xrightarrow{(\mathbf{S}^\top \mathbf{H}(x) \mathbf{S})^\dagger} (\mathbf{S}^\top \mathbf{H}(x) \mathbf{S})^\dagger \mathbf{S}^\top g(x) \xrightarrow{\mathbf{S}} \mathbf{S} (\mathbf{S}^\top \mathbf{H}(x) \mathbf{S})^\dagger \mathbf{S}^\top g(x)$$

Can be computed with directional derivatives:

$$\frac{df(x + \lambda \mathbf{S})}{d\lambda} \Big|_{\lambda=0} = \mathbf{S}^\top g(x) \quad \frac{d^2 f(x + \lambda \mathbf{S})}{d\lambda^2} \Big|_{\lambda=0} = \mathbf{S}^\top \mathbf{H}(x) \mathbf{S}$$

Advantages of RSN:

- Uses second-order information & hence enjoys better dependence on condition number
- Enjoys global convergence theory
- Is a descent method: $f(x_{k+1}) \leq f(x_k)$
- Is a feasible method: $x_k \in \mathcal{Q}$ for all $k \geq 0$
- Applicable for very large d

Example: Single Column Sketches

Let $0 \prec \mathbf{U} \in \mathbb{R}^{d \times d}$ be a symmetric positive definite matrix such that $\mathbf{H}(x) \preceq \mathbf{U}$, $\forall x \in \mathbb{R}^d$. Let $\mathbf{M} = [m_1, \dots, m_d] \in \mathbb{R}^{d \times d}$ be an invertible matrix such that $m_i^\top \mathbf{H}(x) m_i \neq 0$ for all $x \in \mathcal{Q}$ and $i = 1, \dots, d$. If we sample according to

$$\text{Prob}(\mathbf{S}_k = m_i) = p_i := \frac{m_i^\top \mathbf{U} m_i}{\text{Trace}(\mathbf{M}^\top \mathbf{U} \mathbf{M})},$$

then the update on line 5 of Algorithm 1 is given by

$$x_{k+1} = x_k - \frac{1}{\hat{L}} \frac{m_i^\top g(x_k)}{m_i^\top \mathbf{H}(x_k) m_i} m_i, \quad \text{with probability } p_i, \quad (6)$$

costs $\mathcal{O}(d)$ and has linear iteration complexity (10) given by

$$k \geq \max_{x \in \mathcal{Q}} \frac{\text{Trace}(\mathbf{M}^\top \mathbf{U} \mathbf{M})}{\lambda_{\min}^+(\mathbf{H}^{1/2}(x) \mathbf{M} \mathbf{M}^\top \mathbf{H}^{1/2}(x))} \frac{\hat{L}}{\hat{\mu}} \log\left(\frac{1}{\epsilon}\right).$$

6. RSN: Equivalent Viewpoints

1. **Minimization of $T(\cdot, x_k)$ over a random subspace:**

$$x_{k+1} = \arg \min_{x \in \mathbb{R}^d, \lambda \in \mathbb{R}^s} T(x, x_k) \quad (7)$$

subject to $x = x_k + \mathbf{S}_k \lambda$.

2. **Projection of the Newton direction**

$n(x_k) := -\mathbf{H}^\dagger(x_k)g(x_k)$ onto a random subspace:

$$x_{k+1} = \arg \min_{x \in \mathbb{R}^d, \lambda \in \mathbb{R}^s} \left\| x - \left(x_k - \frac{1}{\hat{L}} n(x_k) \right) \right\|_{\mathbf{H}(x_k)}^2 \quad (8)$$

subject to $x = x_k + \mathbf{S}_k \lambda$.

3. **Projection of the current iterate x_k onto a sketched Newton system:**

$$x_{k+1} \in \arg \min_{x \in \mathbb{R}^d} \|x - x_k\|_{\mathbf{H}(x_k)}^2 \quad (9)$$

$$\text{subject to } \mathbf{S}_k^\top \mathbf{H}(x_k)(x - x_k) = -\frac{1}{\hat{L}} \mathbf{S}_k^\top g(x_k).$$

Remark: If $\text{Range}(\mathbf{S}_k) \subset \text{Range}(\mathbf{H}(x_k))$, then x_{k+1} is the unique solution to (9).

7. Convergence Theory

Let $\mathbf{G}(x) := \mathbb{E}_{\mathbf{S} \sim \mathcal{D}} [\mathbf{S} (\mathbf{S}^\top \mathbf{H}(x) \mathbf{S})^\dagger \mathbf{S}]$ and define

$$\rho(x) := \min_{v \in \text{Range}(\mathbf{H}(x))} \frac{\langle \mathbf{H}^{1/2}(x) \mathbf{G}(x) \mathbf{H}^{1/2}(x) v, v \rangle}{\|v\|_2^2}, \quad \rho := \min_{x \in \mathcal{Q}} \rho(x) \leq 1.$$

Global Linear Convergence of RSN

Let $f(x_0) > f_* := \min_x f(x)$. If all assumptions hold, then

$$\mathbb{E}[f(x_k)] - f_* \leq \left(1 - \rho \frac{\hat{\mu}}{\hat{L}}\right)^k (f(x_0) - f_*).$$

Consequently, given $\epsilon > 0$, if $\rho > 0$ then

$$k \geq \frac{1}{\rho \hat{\mu}} \log\left(\frac{1}{\epsilon}\right) \Rightarrow \frac{\mathbb{E}[f(x_k) - f_*]}{f(x_0) - f_*} \leq \epsilon. \quad (10)$$

Sublinear Convergence of RSN

If the assumptions hold with $\hat{L} > \hat{\mu} = 0$ and

$$\mathcal{R} := \inf_{x_* \in \arg \min f} \sup_{x \in \mathcal{Q}} \|x - x_*\|_{\mathbf{H}(x)} < +\infty,$$

and $\rho > 0$ then

$$\mathbb{E}[f(x_k)] - f_* \leq \frac{2\hat{L}\mathcal{R}^2}{\rho k}. \quad (11)$$

Example: RSN includes Newton's method as a special case with $\mathbf{S}_k = \mathbf{I} \in \mathbb{R}^{d \times d}$. In this case, $\rho(x_k) \equiv 1$ and thus (10) recovers the $\frac{\hat{L}}{\hat{\mu}} \log(1/\epsilon)$ complexity given in [1] and (11) gives a new sublinear result.

Sufficient Condition for $\rho > 0$

If (5) holds and $\text{Range}(\mathbf{H}(x_k)) \subset \text{Range}(\mathbb{E}[\mathbf{S}_k \mathbf{S}_k^\top])$, then $\rho > 0$, and $\rho = \lambda_{\min}^+(\mathbb{E}_{\mathbf{S} \sim \mathcal{D}} [\mathbf{H}^{1/2}(x_k) \mathbf{S}_k (\mathbf{S}_k^\top \mathbf{H}(x_k) \mathbf{S}_k)^\dagger \mathbf{S}_k^\top \mathbf{H}^{1/2}(x_k)])$

Example: Generalized Linear Models

Let $0 \leq u \leq \ell$. Let $\phi_i: \mathbb{R} \mapsto \mathbb{R}_+$ be a twice differentiable function such that

$$u \leq \phi_i''(t) \leq \ell, \quad \text{for } i = 1, \dots, n. \quad (12)$$

Let $a_i \in \mathbb{R}^d$ for $i = 1, \dots, n$ and $\mathbf{A} = [a_1, \dots, a_n] \in \mathbb{R}^{d \times n}$. We say that $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a generalized linear model when

$$f(x) = \frac{1}{n} \sum_{i=1}^n \phi(a_i^\top x) + \frac{\lambda}{2} \|x\|_2^2. \quad (13)$$

f is \hat{L} -smooth and $\hat{\mu}$ -convex relative to its Hessian with

$$\hat{L} = \frac{\ell \sigma_{\max}^2(\mathbf{A}) + n\lambda}{u \sigma_{\max}^2(\mathbf{A}) + n\lambda} \quad \text{and} \quad \hat{\mu} = \frac{u \sigma_{\max}^2(\mathbf{A}) + n\lambda}{\ell \sigma_{\max}^2(\mathbf{A}) + n\lambda}. \quad (14)$$

RSN has iteration complexity (10) given by

$$k \geq \frac{1}{\rho} \left(\frac{\ell \sigma_{\max}^2(\mathbf{A}) + n\lambda}{u \sigma_{\max}^2(\mathbf{A}) + n\lambda} \right)^2 \log\left(\frac{1}{\epsilon}\right). \quad (15)$$

8. Experiments

We compare RSN to Gradient descent (GD), accelerated gradient descent (AGD) [2] and full Newton method. For RSN we use coordinate sketches defined by $\mathbf{S}_k \in \{0, 1\}^{d \times s}$, with exactly one non-zero entry per row and per column of \mathbf{S}_k .

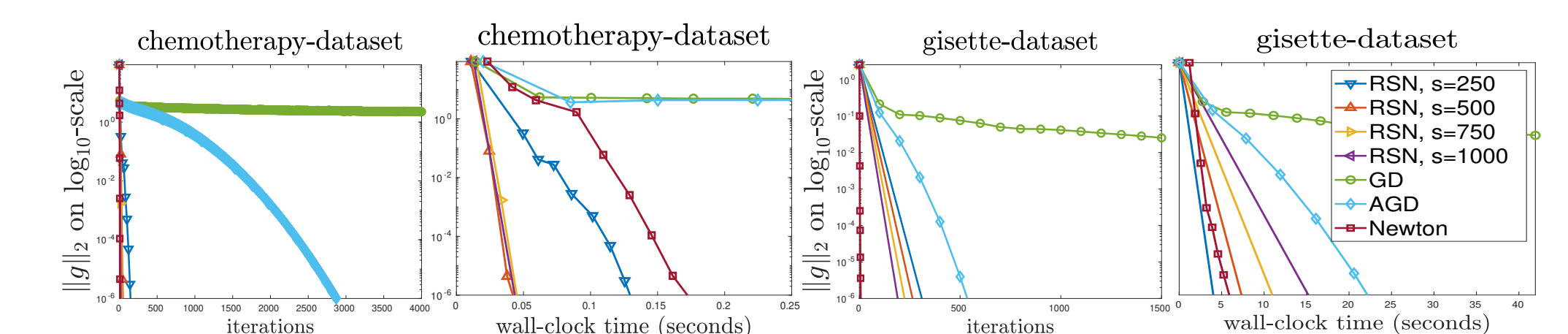


Figure: Highly dense problems, favoring RSN methods.

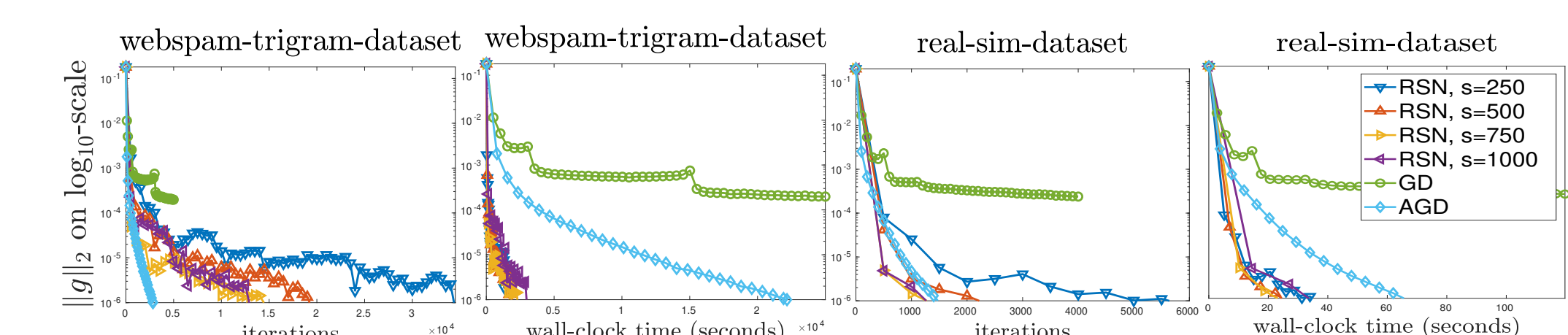


Figure: Moderately sparse problems favor the RSN method. The full Newton method is infeasible due to high dimensionality.

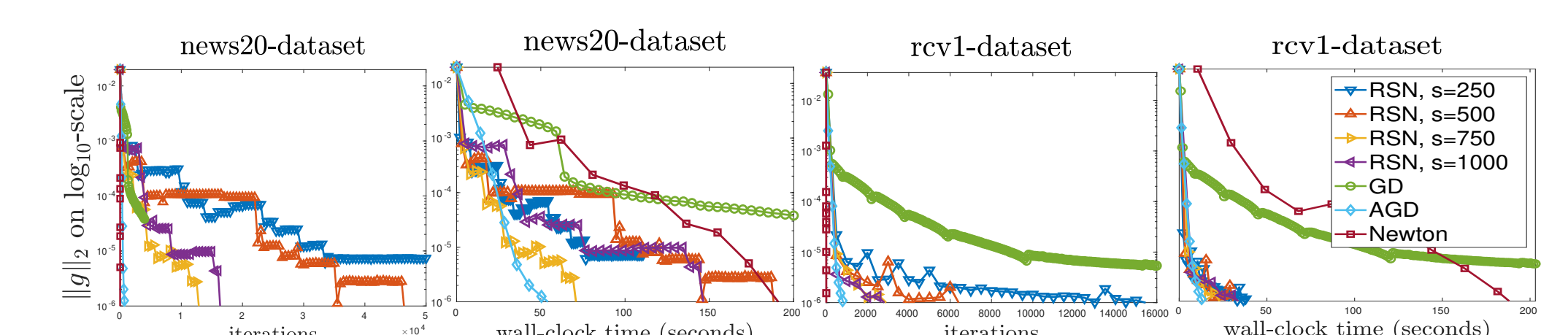


Figure: Due to extreme sparsity, accelerated gradient is competitive with the Newton type methods.

References

- [1] S. P. Karimireddy, S. U. Stich, and M. Jaggi. Global linear convergence of Newton's method without strong-convexity or Lipschitz gradients. *arXiv:1806.0041*, 2018.
- [2] Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Springer Publishing Company, Incorporated, 1 edition, 2014.