



西安电子科技大学
XIDIAN UNIVERSITY



西安交通大学
XI'AN JIAOTONG UNIVERSITY

UNIVERSITÉ
SORBONNE
PARIS NORD



澳門科技大學
MACAU UNIVERSITY OF SCIENCE AND TECHNOLOGY



THE HONG KONG
POLYTECHNIC UNIVERSITY
香港理工大學

Globe Q-linear Gauss-Newton Method for Overparameterized Non-convex Matrix Sensing

Xixi Jia¹, Fangchen Feng², Deyu Meng^{3,4}, Defeng Sun⁵

hsijixidian@gmail.com

¹ School of Mathematics and Statistics, Xidian University;

² L2TI Laboratory, University Sorbonne Paris Nord;

³ School of Mathematics and Statistics, Xi'an Jiaotong University;

⁴ Macao Institute of Systems Engineering, Macau University of Science and Technology

⁵ Department of Applied Mathematics, The Hong Kong Polytechnic University

Review of the Overparameterized Non-convex Matrix Sensing

Matrix sensing aims to recover an unknown low-rank matrix $M \in \mathbb{R}^{n \times n}$ from its linear measurement $\mathbf{b} = \mathcal{A}(M)$ by solving the following optimization problem

$$\min_{U \in \mathbb{R}^{n \times d}, V \in \mathbb{R}^{n \times d}} f(U, V) := \frac{1}{2} \|\mathcal{A}(UV^\top) - \mathbf{b}\|_F^2,$$

where $\text{rank}(M) = r \ll n$, of particular interest is the case where $d > r$.

Review of the Overparameterized Non-convex Matrix Sensing

Matrix sensing aims to recover an unknown low-rank matrix $M \in \mathbb{R}^{n \times n}$ from its linear measurement $\mathbf{b} = \mathcal{A}(M)$ by solving the following optimization problem

$$\min_{U \in \mathbb{R}^{n \times d}, V \in \mathbb{R}^{n \times d}} f(U, V) := \frac{1}{2} \|\mathcal{A}(UV^\top) - \mathbf{b}\|_F^2,$$

where $\text{rank}(M) = r \ll n$, of particular interest is the case where $d > r$. **The problem is challenging due to the following reasons:**

- ① The optimization problem is non-convex and non-smooth;
- ② The saddle points can slow down the converges of the gradient based algorithms;
- ③ Overparameterization can further degenerate the convergence of GD from linear rate to sub-linear rate;

Review of the Overparameterized Non-convex Matrix Sensing

The results in Table 1 and Figure 3 illustrate the performance of current gradient-based methods on this canonical problem.

Table 1: Comparisons of iteration complexity, with κ as the condition number of the $n \times n$ matrix. “init.” denotes initialization.

Algorithm	init.	iteration complexity
GD [20]	random	$\kappa^{11} \log(\kappa^2/n) + \kappa^{10} \log(\kappa^6/\varepsilon)$
PrecGD [15]	spectral	$\log(1/\varepsilon)$
ScaledGD(λ)[21]	random	$\log \kappa \cdot \log(\kappa n) + \log(1/\varepsilon)$
AGN	random	$\log(1/\varepsilon)$

Review of the Overparameterized Non-convex Matrix Sensing

The results in Table 1 and figure illustrate the performance of current gradient-based methods on this canonical problem.

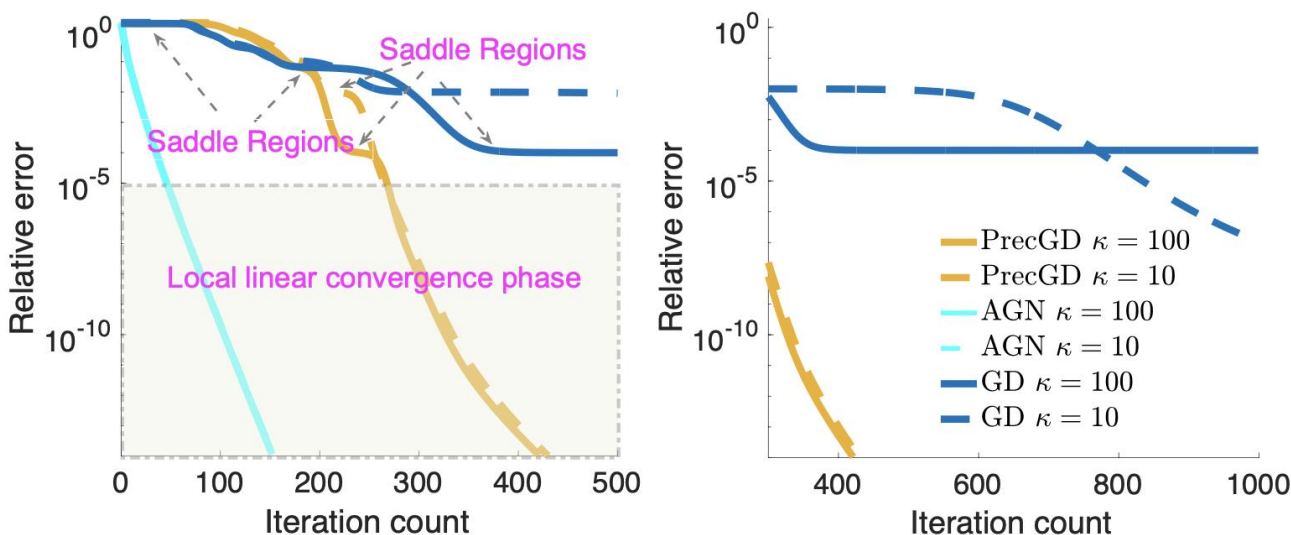


Figure 1: Comparison of convergence for PrecGD, GD, and AGN across various condition numbers, with the right subfigure extending the left by iterating from 300 to 1000.

Approximated Gauss-Newton method (AGN)

2.1 AGN for asymmetric matrix sensing

We unify the formulation of symmetric and asymmetric low rank matrix sensing (LRMS) into a single, simplified expression:

$$\min_{X \in \mathbb{R}^{2n \times d}} \psi(X) := \frac{1}{2} \|\mathcal{A}(PXX^\top Q) - \mathbf{b}\|_2^2$$

where $P = \begin{bmatrix} I & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n \times 2n}$, $Q = \begin{bmatrix} \mathbf{0} \\ I \end{bmatrix} \in \mathbb{R}^{2n \times n}$ and $I \in \mathbb{R}^{n \times n}$ is the identity matrix. The case $X = \begin{bmatrix} U \\ V \end{bmatrix}$ with $U \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{n \times n}$ corresponds to the asymmetric matrix sensing.

Approximated Gauss-Newton method (AGN)

2.1 AGN for asymmetric matrix sensing

We unify the formulation of symmetric and asymmetric low rank matrix sensing (LRMS) into a single, simplified expression:

$$\min_{X \in \mathbb{R}^{2n \times d}} \psi(X) := \frac{1}{2} \|\mathcal{A}(PXX^\top Q) - \mathbf{b}\|_2^2$$

By employing the Gauss-Newton method, one can update the variable as $X_{t+1} = X_t + \eta \Delta(X_t)$ where

$$\Delta(X_t) = \arg \min_{\Delta \in \mathbb{R}^{2n \times d}} \frac{1}{2} \|\mathcal{B}(\Delta, X_t) + \mathcal{B}(X_t, \Delta) + \mathcal{B}(X_t, X_t) - \mathbf{b}\|_2^2$$

Approximated Gauss-Newton method (AGN)

2.1 AGN for asymmetric matrix sensing

By employing the Gauss-Newton method, one can update the variable

as $X_{t+1} = X_t + \eta \Delta(X_t)$ where

$$\Delta(X_t) = \arg \min_{\Delta \in \mathbb{R}^{2n \times d}} \frac{1}{2} \|\mathcal{B}(\Delta, X_t) + \mathcal{B}(X_t, \Delta) + \mathcal{B}(X_t, X_t) - \mathbf{b}\|_2^2$$

We apply a Gauss-Seidel method to solve the above least square problem as

$$X_{t+\frac{1}{2}} = X_t + \eta \Delta(X_t), \quad \Delta(X_t) = \arg \min_{\Delta \in \mathbb{R}^{2n \times d}} \frac{1}{2} \|\mathcal{B}(\Delta, X_t) + \mathcal{B}(X_t, X_t) - \mathbf{b}\|_2^2,$$

$$X_{t+1} = X_{t+\frac{1}{2}} + \eta \Delta(X_{t+\frac{1}{2}}), \quad \Delta(X_{t+\frac{1}{2}}) = \arg \min_{\Delta \in \mathbb{R}^{2n \times d}} \frac{1}{2} \|\mathcal{B}(X_{t+\frac{1}{2}}, \Delta) + \mathcal{B}(X_{t+\frac{1}{2}}, X_{t+\frac{1}{2}}) - \mathbf{b}\|_2^2.$$

Approximated Gauss-Newton method (AGN)

2.1 AGN for asymmetric matrix sensing

Thanks to the RIP condition of the LRMS problem, one can approximate the Gauss-Newton direction by the following

$$\hat{\Delta}(X_t) = \arg^+ \min_{\Delta \in \mathbb{R}^{2n \times d}} \frac{1}{2} \|\hat{\mathcal{B}}(\Delta, X_t) - \mathcal{A}^*(\mathcal{B}(X_t, X_t) - \mathbf{b})\|_F^2,$$

$$\hat{\Delta}(X_{t+\frac{1}{2}}) = \arg^+ \min_{\Delta \in \mathbb{R}^{2n \times d}} \frac{1}{2} \|\hat{\mathcal{B}}(X_{t+\frac{1}{2}}, \Delta) - \mathcal{A}^*(\mathcal{B}(X_{t+\frac{1}{2}}, X_{t+\frac{1}{2}}) - \mathbf{b})\|_F^2,$$

where $\hat{\mathcal{B}}(\Delta, X_t) = P\Delta X_t^\top Q$ and $\hat{\mathcal{B}}(X_{t+\frac{1}{2}}, \Delta) = PX_{t+\frac{1}{2}}\Delta^\top Q$. \arg^+ denotes the minimum norm solution.

Approximated Gauss-Newton method (AGN)

2.1 AGN for asymmetric matrix sensing

Lemma 1. (Descent lemma) *For asymmetric matrix sensing, as long as $0 < \eta \leq 2/(1 + \delta)$ and the Assumption 1 is satisfied, then there exists positive constant $\ell = (2\eta - (1 + \delta)\eta^2)/2$ such that*

$$\psi(X_{t+\frac{1}{2}}) \leq \psi(X_t) - \ell \|\hat{\mathcal{B}}(\hat{\Delta}(X_t), X_t)\|_F^2,$$

$$\psi(X_{t+1}) \leq \psi(X_{t+\frac{1}{2}}) - \ell \|\hat{\mathcal{B}}(X_{t+\frac{1}{2}}, \hat{\Delta}(X_{t+\frac{1}{2}}))\|_F^2.$$

Lemma 1 suggests that AGN with a constant step-size is indeed a descent method for the overparameterized LRMS.

Approximated Gauss-Newton method(AGN)

2.2 AGN for symmetric matrix sensing

The case $X = \begin{bmatrix} U \\ U \end{bmatrix}$ corresponds to the symmetric matrix sensing.

There are two different ways to deal with symmetric MS, by symmetric parameterization $X \in \mathcal{C}$ or asymmetric parameterization $X \in \mathbb{R}^{2n \times d}$,

where

$$\mathcal{C} = \left\{ Z \mid Z = \begin{bmatrix} U \\ U \end{bmatrix}, U \in \mathbb{R}^{n \times d} \right\} \subset \mathbb{R}^{2n \times d}$$

and the update is given by

$$\tilde{\Delta}(X_t) = \arg \min_{\Delta \in \mathcal{C}} \frac{1}{2} \|\hat{\mathcal{B}}(\Delta, X_t) - \mathcal{A}^*(\mathcal{B}(X_t, X_t) - \mathbf{b})\|_F^2$$

Approximated Gauss-Newton method(AGN)

2.2 AGN for symmetric matrix sensing

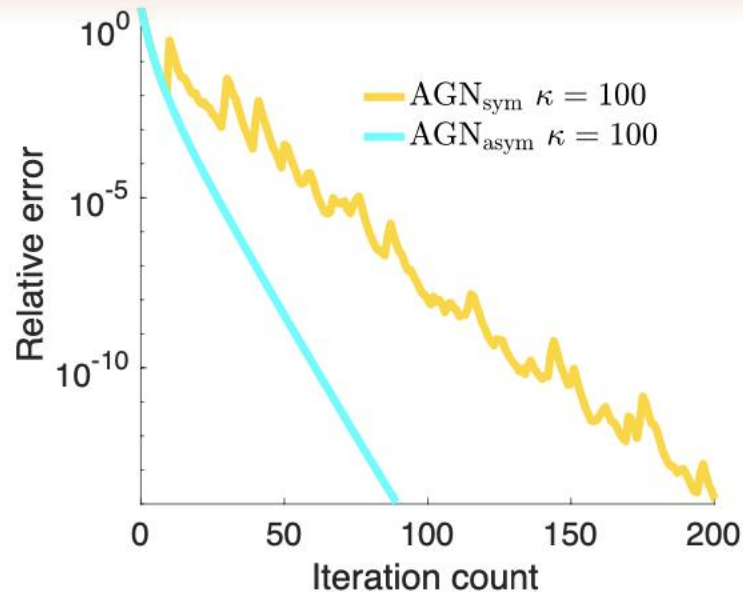


Figure 2. Convergence of AGN under Sym. and Asym. parameterization of symmetric LRMS.

Saddle point analysis

We consider the population risk of the LRMS problem as solving the following problem

$$\min_{X \in \mathbb{R}^{2n \times d}} \frac{1}{2} \|PXX^\top Q - M\|_2^2.$$

The objective function corresponds to $g(U, V) = \frac{1}{2} \|UV^\top - M\|_F^2$ where $U \in \mathbb{R}^{n \times d}$ and $V \in \mathbb{R}^{n \times d}$. The saddle point of the non-convex objective is denoted by $(U_s, V_s) \in \mathcal{S}$ where \mathcal{S} is defined as

$$\mathcal{S} = \{(U_s, V_s) | U_s V_s^\top = \Phi \mathcal{M}(\Sigma) \Psi^\top, M = \Phi \Sigma \Psi^\top, \mathcal{M} \in \mathfrak{M}/\mathfrak{J}\}$$

where $M = \Phi \Sigma \Psi^\top$ is the SVD of the matrix M .

Saddle point analysis

We now establish the following theorem to characterize the saddle point analysis.

Theorem 1. *Assume that M is rank-1, the point (\hat{U}, \hat{V}) with $\hat{U} = U_s + \varepsilon N_u, \hat{V} = V_s + \varepsilon N_v$ is at the vicinity of the saddle point $(U_s, V_s) \in \mathcal{S}$ and ε is sufficiently small, N_u and N_v are random Gaussian matrices that follow a standard normal distribution. Then with high probability we have the following results*

$$\text{(GD)} \quad \|\nabla g\|_F^2 = o(\varepsilon)e_s + o(\varepsilon^2),$$

$$\text{(AGN)} \quad \begin{cases} \|\nabla g_{\hat{U}}(\hat{V}^\top \hat{V})^{-\frac{1}{2}}\|_F^2 = \Theta(1)e_s + o(\varepsilon^2), \\ \|\nabla g_{\hat{V}}(\hat{U}^\top \hat{U})^{-\frac{1}{2}}\|_F^2 = \Theta(1)e_s + o(\varepsilon^2), \end{cases}$$

where $e_s = \|U_s V_s^\top - M\|_F^2$. Furthermore, by constraining M to be positive semi-definite and $\hat{U} = \hat{V}, U_s = V_s$ (for symmetric matrix sensing), for bounded constant $c > 0$, we have

$$\text{(ScaledGD}(\lambda)\text{)} \quad \|\nabla g_{\hat{U}}(\hat{U}^\top \hat{U} + \lambda I)^{-\frac{1}{2}}\|_F^2 = \Theta\left(\frac{\varepsilon^2}{\varepsilon^2 + \lambda/c}\right) e_s + o(\varepsilon^2).$$

Saddle point analysis

We now establish the following theorem to characterize the saddle point analysis.

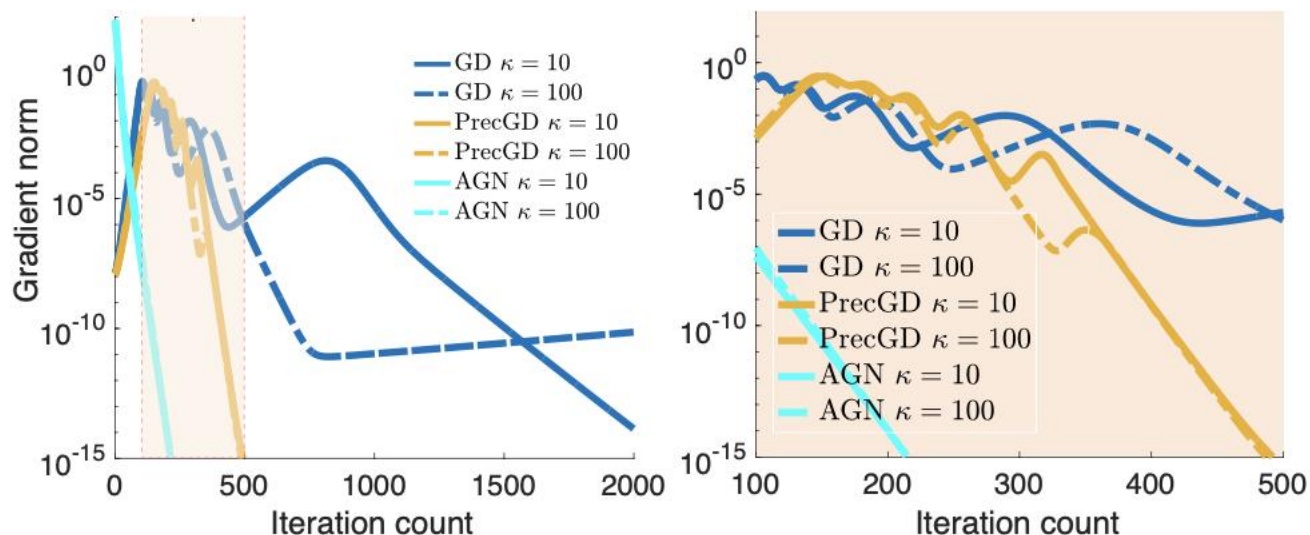


Figure 3: Illustration of the gradient norm for GD, PrecGD, and the proposed AGN, with the right subfigure showing a zoomed-in region of the left for iterations from 100 to 500.

Main results

We now present the convergence result of the proposed AGN method.

Theorem 2 (Global Q-linear convergence). *Under the Assumption 1 and Assumption 2. Let ψ^* be the global minimal value of $\psi(X)$ in Eq. (6) and $X_t, \forall t > 0$ is generated by Algorithm 1, then there exists constants $1 \geq \tau > 0$ such that*

$$\psi(X_{t+1}) - \psi^* \leq c_q[\psi(X_t) - \psi^*], \forall t > 0, \quad (24)$$

where $c_q = (1 - \hat{\ell} \frac{1-\delta}{1+\delta} \tau) < 1$ and $\hat{\ell} = 2\eta - (1 + \delta)\eta^2$. Meanwhile, if $\delta = 0, \eta = 1$, c_q becomes 0.

Main results

We now present the convergence result of the proposed AGN method.

Theorem 2 (Global Q-linear convergence). *Under the Assumption 1 and Assumption 2. Let ψ^* be the global minimal value of $\psi(X)$ in Eq. (6) and $X_t, \forall t > 0$ is generated by Algorithm 1, then there exists constants $1 \geq \tau > 0$ such that*

$$\psi(X_{t+1}) - \psi^* \leq c_q[\psi(X_t) - \psi^*], \forall t > 0, \quad (24)$$

where $c_q = (1 - \hat{\ell} \frac{1-\delta}{1+\delta} \tau) < 1$ and $\hat{\ell} = 2\eta - (1 + \delta)\eta^2$. Meanwhile, if $\delta = 0, \eta = 1$, c_q becomes 0.

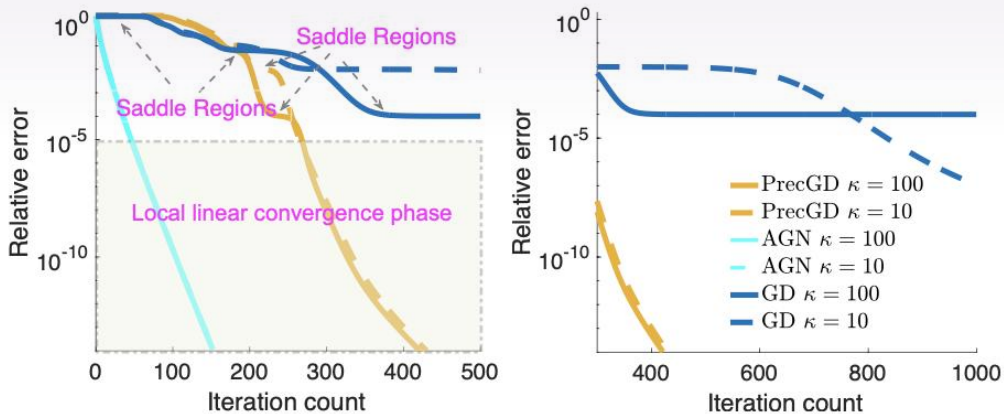


Figure 1: Comparison of convergence for PrecGD, GD, and AGN across various condition numbers, with the right subfigure extending the left by iterating from 300 to 1000.



THE HONG KONG
POLYTECHNIC UNIVERSITY
香港理工大學



UNIVERSITÉ
SORBONNE
PARIS NORD

Thank you

Globe Q-linear Gauss-Newton Method for Overparameterized Non-convex Matrix Sensing