

Contents

1	Introduction	1
2	The Anderson mixing scheme	3
3	The Min-AM methods	3
3.1	The basic Min-AM	3
3.2	The restarted Min-AM	4
3.3	Eigenvalue estimates and the choice of mixing parameter	5
3.4	The stochastic Min-AM	6
4	Experiments	8
5	Conclusion	10
A	Comparisons with conventional methods	16
A.1	Relationship between AM-I, Min-AM, and Newton’s method	16
A.2	Comparison with the conjugate gradient/residual method	16
A.3	Comparison with the BFGS method	18
A.4	Comparison with the momentum-based method	18
A.5	Comparison with related variants of Anderson mixing	19
B	Details of the basic Min-AM	19
B.1	Derivation of the basic Min-AM	19
B.2	Proof of Theorem 1	21
B.3	Convergence analysis	24
C	Details of the restarted Min-AM	25
C.1	Proof of Theorem 2	25
D	Details of the eigenvalue estimation procedure	33
D.1	Analysis of the quadratic case	33
D.2	Proof of Theorem 3	34
E	Details of the stochastic Min-AM	36
E.1	Some useful results	37
E.2	Proof of Theorem 4	40
E.3	Proof of Theorem 5	41
F	Experimental details	42
F.1	Strongly convex quadratic problem	42
F.2	Regularized logistic regression problem	44
F.3	Deep neural network training problem	46

F.3.1	Hyperparameter setting of the stochastic Min-AM	46
F.3.2	Experiments on CIFAR-10 and CIFAR-100	46
F.3.3	Experiment of training ResNet50 on ImageNet	51
F.3.4	An additional experiment: adversarial training	51

G Limitations

52

A Comparisons with conventional methods

We compare the proposed Min-AM method with other related methods, including Newton’s method, the (nonlinear) conjugate gradient method, BFGS method, momentum-based method, and some other variants of Anderson mixing.

A.1 Relationship between AM-I, Min-AM, and Newton’s method

We first give a new interpretation of the AM-I method here, which reveals the relationship between AM-I and Newton’s method.

Recall that the historical information is stored in $X_k, R_k \in \mathbb{R}^{d \times m}$. We have the decomposition of the solution space by $\mathbb{R}^d = \text{range}(X_k) \oplus \text{range}(X_k)^\perp$. Define $V_k \in \mathbb{R}^{d \times (d-m)}$ whose columns form an orthonormal basis of $\text{range}(X_k)^\perp$. Then for any $x \in \mathbb{R}^d$, we have $x = x_k - X_k \gamma - V_k \eta$, where $\gamma \in \mathbb{R}^m, \eta \in \mathbb{R}^{d-m}$. As a result,

$$\min_{x \in \mathbb{R}^d} f(x) = \min_{\gamma \in \mathbb{R}^m, \eta \in \mathbb{R}^{d-m}} f(x_k - X_k \gamma - V_k \eta). \quad (29)$$

Suppose that f is twice differentiable and $X_k^T \nabla^2 f(x_k) X_k$ is nonsingular. We apply Newton’s method in the low dimensional subspace $\text{range}(X_k)$ followed by a gradient descent with stepsize β_k in the complementary subspace $\text{range}(X_k)^\perp$. Then

$$\bar{x}_k = x_k - X_k \gamma_k, \quad x_{k+1} = \bar{x}_k - V_k \eta_k, \quad (30)$$

where $\gamma_k = (X_k^T \nabla^2 f(x_k) X_k)^{-1} X_k^T \nabla f(x_k)$, and $\eta_k = \beta_k V_k^T \nabla f(\bar{x}_k)$. When $m \ll d$, the scheme (30) is more economical than applying the Newton’s method in the whole solution space \mathbb{R}^d .

Now, consider the simple case that f is a quadratic function: $f(x) = \frac{1}{2} x^T A x - b^T x$. The residual $r_k = -\nabla f(x_k) = b - A x_k$. Then $r_k - r_{k-1} = -A(x_k - x_{k-1})$, which implies $R_k = -A X_k$. It follows that $\gamma_k = (X_k^T R_k)^{-1} X_k^T r_k$. Also, $\bar{r}_k := r_k - R_k \gamma_k = -\nabla f(\bar{x}_k)$. So $\bar{r}_k \perp \text{range}(X_k)$ due to the choice of γ_k , i.e., a kind of the Galerkin’s projection condition [53]. Let U_k be the matrix whose columns form an orthonormal basis of $\text{range}(X_k)$. Then $V_k V_k^T = I - U_k U_k^T$ and $U_k^T \bar{r}_k = 0$. For the gradient descent step, we have $\eta_k = -\beta_k V_k^T \bar{r}_k$, which implies

$$x_{k+1} = \bar{x}_k - V_k \eta_k = \bar{x}_k + \beta_k V_k V_k^T \bar{r}_k = \bar{x}_k + \beta_k (I - U_k U_k^T) \bar{r}_k = \bar{x}_k + \beta_k \bar{r}_k,$$

namely a mixing step. Therefore, the AM-I method coincides with the scheme (30) for minimizing a quadratic function. Since Min-AM is essentially equivalent to the full-memory AM-I when minimizing a strongly convex quadratic function, we know Min-AM is also closely related to the scheme (30). In the general case, if the objective function f can be well approximated by a quadratic function in a local region around the optima, it is expected that AM-I and Min-AM can behave similarly to the scheme (30) when the iterates enter this local region.

A.2 Comparison with the conjugate gradient/residual method

The conjugate gradient (CG) method [32] and the conjugate residual (CR) method [53, Algorithm 6.20] are classical methods for solving SPD linear systems, where the CG method is extended to solving unconstrained optimization [45]. We first discuss the connection between the proposed Min-AM method and the CG/CR method.

Connection between Min-AM and CG/CR. For solving linear systems, it has been established in [60] that the full-memory AM-I and AM are essentially equivalent to the full orthogonalization

method (FOM) [52] and GMRES [51], respectively. If the linear system is SPD, both FOM and GMRES methods can be simplified to have short-term recurrences: FOM is equivalent to CG, and GMRES is equivalent to CR. Thus, the full-memory AM-I and AM are essentially equivalent to CG and CR respectively, for solving SPD linear systems. However, CG and CR are much more efficient in terms of memory and per-step computational cost, which motivates the development of the proposed Min-AM method that also has short-term recurrences. From Theorem 1, for solving SPD linear systems, Min-AM is essentially equivalent to CG, FOM, and the full-memory AM-I. Compared to CG, Min-AM has the advantage that it does not need explicit matrix-vector products to determine the step size. The update scheme of Min-AM only depends on the historical iterations, which makes it easier to extend to the nonlinear case. Min-AM can also have a similar convergence behaviour to CR for solving SPD linear systems, due to the close relationship between FOM and GMRES [53, Section 6.5.7].

Now, we consider the unconstrained optimization problem

$$\min_{x \in \mathbb{R}^d} f(x),$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is continuously differentiable and bounded from below. The nonlinear CG (NCG) generates $\{x_k\}$ by the update scheme

$$x_{k+1} = x_k + \alpha_k p_k, \quad (31)$$

where the step size α_k is usually obtained by a line search, and the searching direction p_k is constructed by

$$p_{k+1} = r_k + \beta_k p_k, \quad p_0 = r_0. \quad (32)$$

Here, $r_k := -\nabla f(x_k)$ and the β_k is the momentum term. There are several choices of β_k , for example, the Fletcher-Reeves variant $\beta_k = \frac{\|r_{k+1}\|_2^2}{\|r_k\|_2^2}$ and the Polak-Ribière variant $\beta_k = \frac{r_{k+1}^T (r_{k+1} - r_k)}{\|r_k\|_2^2}$. (See [45].) However, the step size α_k is not easy to obtain. If f is a quadratic function, let A be the Hessian of f . The computation of α_k involves a matrix-vector product, e.g.,

$$\alpha_k = \frac{r_k^T p_k}{p_k^T A p_k},$$

which ensures $r_{k+1} \perp p_k$. If finite difference is used to compute $A p_k$, it requires two gradient evaluations. If f is a general nonlinear function, a line search is necessary to ensure the convergence of NCG: α_k is chosen as an approximate solution to the problem

$$\min_{\alpha \geq 0} f(x_k + \alpha p_k), \quad (33)$$

where p_k is supposed to be a descent direction. The d -step quadratic convergence proved in [17], which is similar to our result of the restarted Min-AM, was only established under the assumption of an exact line search, namely, (33) is exactly solved. More practical choices of α_k require some conditions to be satisfied. For example, the strong Wolfe conditions are

$$\begin{aligned} f(x_k + \alpha_k p_k) &\leq f(x_k) + c_1 \alpha_k \nabla f(x_k)^T p_k, \\ |\nabla f(x_k + \alpha_k p_k)^T p_k| &\leq c_2 |\nabla f(x_k)^T p_k|, \end{aligned}$$

where $0 < c_1 < c_2 < 1$. Since the line search can incur many times of the backtracking procedure, the total number of gradient evaluations for NCG can be large.

In contrast, our proposed Min-AM method does not need explicit Hessian-vector products or line search to determine the step size, while the convergence results (cf. Theorem 2) match those of NCG with exact line search. Therefore, Min-AM is more economical than NCG, and can have better convergence than the NCG with inexact line search. Also, it should be pointed out that even in strongly convex quadratic optimization, Min-AM generally does not generate the same sequence of iterates as CG: They are equivalent in the sense that $x_k^{(1)} = x_k^{\text{CG}}$, where $x_k^{(1)}$ is the intermediate step in (8a).

A.3 Comparison with the BFGS method

The BFGS constructs the approximate inverse Hessian by solving

$$H_k = \arg \min_H \|H - H_{k-1}\|_{F(W)}, \quad \text{s.t. } Hy_k = s_k, \quad H = H^T, \quad (34)$$

where $s_k := x_k - x_{k-1}$, $y_k := \nabla f(x_k) - \nabla f(x_{k-1})$. The norm $\|\cdot\|_{F(W)}$ is the weighted Frobenius norm (i.e., $\|X\|_{F(W)}^2 := \|W^{1/2} X W^{1/2}\|_F^2$ for a matrix $X \in \mathbb{R}^{d \times d}$) and the weight matrix W satisfies $W s_k = y_k$. The solution is

$$H_k = \left(I - \frac{s_k y_k^T}{y_k^T s_k} \right) H_{k-1} \left(I - \frac{y_k s_k^T}{y_k^T s_k} \right) + \frac{s_k s_k^T}{y_k^T s_k}. \quad (35)$$

Note that the $H_k \in \mathbb{R}^{d \times d}$ is recursively constructed and can be dense during the later iterations. To reduce memory overhead, the limited-memory BFGS [39] is often used. In memoryless BFGS, the previous approximate Hessian is always reset to I_d , i.e.,

$$H_k = \left(I - \frac{s_k y_k^T}{y_k^T s_k} \right) \left(I - \frac{y_k s_k^T}{y_k^T s_k} \right) + \frac{s_k s_k^T}{y_k^T s_k}. \quad (36)$$

Memoryless BFGS has minimal memory size in L-BFGS, but has no equivalence to BFGS. It is pointed in [45] that for the update $x_{k+1} = x_k - \alpha_k H_k \nabla f(x_k)$ of memoryless BFGS, if α_k is chosen by an exact line search, then memoryless BFGS is equivalent to CG for quadratic function. However, an exact line search can incur prohibitive cost in practice.

The Min-AM method has a similar form to that of memoryless BFGS, but uses a recursively constructed vector pair. Min-AM is essentially equivalent to the full-memory AM-I and CG for minimizing strongly convex quadratic functions. In fact, let $P_k = (p_1, \dots, p_k)$, $Q_k = (q_1, \dots, q_k)$, and define

$$H_k^A = -P_k (P_k^T Q_k)^{-1} P_k^T + \beta_k (I - P_k (P_k^T Q_k)^{-1} Q_k^T) (I - Q_k (P_k^T Q_k)^{-1} P_k^T), \quad (37)$$

assuming $P_k^T Q_k$ is nonsingular. We call the iterations defined by $x_{k+1} = x_k + H_k^A r_k$ as Scheme A. Clearly, $H_k^A Q_k = -P_k$ holds. Moreover, it can be proved by direct computation (using the properties in Theorem 1) that for strongly convex quadratic optimization, if fixed β_k is used, i.e., $\beta_k \equiv \beta$, where β is a constant, we have

$$H_k^A = -\frac{p_k p_k^T}{p_k^T q_k} + \left(I - \frac{p_k q_k^T}{p_k^T q_k} \right) H_{k-1}^A \left(I - \frac{q_k p_k^T}{p_k^T q_k} \right), \quad (38)$$

starting from $H_0^A := \beta I$. Hence, the H_k^A in (37) solves

$$H_k^A = \arg \min_H \|H - H_{k-1}^A\|_{F(W)}, \quad \text{s.t. } H q_k = -p_k, \quad H = H^T, \quad (39)$$

where the weight matrix satisfies $W p_k = -q_k$. The (39) is similar to (34), but based on modified vector pairs. (Note that we use $r_k - r_{k-1} = -(\nabla f(x_k) - \nabla f(x_{k-1}))$ to construct q_k , so there is a difference of sign.) It can also be verified that for strongly convex quadratic optimization, $H_k r_k = H_k^A r_k$, where H_k is defined in (10). So the iterations of Min-AM are identical to those of Scheme A in this case. In this sense, though only using one vector pair, Min-AM implicitly constructs H_k^A which can well approximate the inverse Hessian.

A.4 Comparison with the momentum-based method

The momentum-based method is similar to the scheme defined in (31) and (32), but the momentum β_k is different from that in CG. We take the Nesterov's accelerated gradient (NAG) method as an example. The update scheme of NAG is defined as

$$y_{k+1} = x_k - \frac{1}{L} \nabla f(x_k), \quad (40a)$$

$$x_{k+1} = y_{k+1} + \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} (y_{k+1} - y_k), \quad (40b)$$

where L and μ are the Lipschitz constant and the strong convexity constant of ∇f , respectively. Unlike Min-AM, NAG generally does not form a symmetric approximation to the Hessian. In Figure 3, we show the convergence behaviours of NAG with different settings of μ and L . It is found that the setting of L and μ can have a large effect on the convergence. On the contrary, Min-AM forms symmetric Hessian approximations and can give useful information of L and μ via an economical eigenvalue estimation procedure.

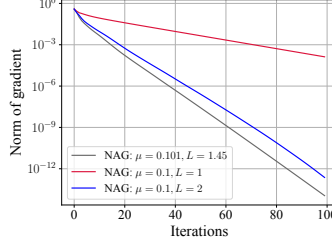


Figure 3: Convergence behaviours of NAG with different settings of L and μ , for the regularized logistic regression on the madelon dataset.

A.5 Comparison with related variants of Anderson mixing

Compared with the original AM, AM-I, the regularized nonlinear acceleration scheme [55], and the stochastic AM [62], the proposed Min-AM has the minimal memory size, the same as that of AM(1). Nonetheless, Min-AM incorporates more historical information than AM(1) and AM-I(1), and is equivalent to the full-memory AM-I in solving strongly convex quadratic optimization. The convergence analysis in Section 3.4 also shows that the stochastic Min-AM has similar convergence to that of the stochastic AM in theory.

The short-term recurrence AM (ST-AM) [63] stores two vector pairs and is equivalent to the full-memory AM in solving strongly convex quadratic optimization, but the approximated Hessian from ST-AM is generally not symmetric. It is also worth pointing out that for AM and ST-AM, their existing theoretical results do not show the significantly better convergence rate than the simple gradient descent method for unconstrained optimization.

Our proposed Min-AM further reduces the memory size of ST-AM and the faster convergence than gradient descent for general nonlinear optimization is rigorously justified in theory.

B Details of the basic Min-AM

The procedure of the basic Min-AM is described in Algorithm 2. Next, we give the derivation of the basic Min-AM, and prove the related theoretical properties.

B.1 Derivation of the basic Min-AM

We give the derivation of the update scheme (10) for solving the quadratic optimization $\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{2}x^T Ax - b^T x$ here. Recall that the steps in one iteration of the basic Min-AM are:

$$x_k^{(1)} = x_k - p_k \Gamma_k^{(1)}, \quad x_k^{(2)} = x_k^{(1)} + \beta_k r_k^{(1)}, \quad x_{k+1} = x_k^{(2)} - p_k \Gamma_k^{(2)}, \quad (41)$$

where $r_k^{(1)} = r_k - q_k \Gamma_k^{(1)}$ and $\beta_k > 0$. Define $r_k^{(2)} = r_k^{(1)} - \beta_k A r_k^{(1)}$ and $r_k^{(3)} = r_k^{(2)} - q_k \Gamma_k^{(2)}$. The $\Gamma_k^{(1)}$ and $\Gamma_k^{(2)}$ are determined by imposing the projection conditions:

$$r_k^{(1)} \perp p_k, \quad r_k^{(3)} \perp p_k.$$

Here, we define $p_1 = \Delta x_0$, $q_1 = \Delta r_0$, and for $k \geq 2$, the construction of p_k, q_k at the beginning of the k -th iteration is

$$p_k = \Delta x_{k-1} - p_{k-1} \zeta_k, \quad q_k = \Delta r_{k-1} - q_{k-1} \zeta_k,$$

where $\zeta_k = (p_{k-1}^T q_{k-1})^{-1} p_{k-1}^T \Delta r_{k-1}$, assuming $p_{k-1}^T q_{k-1} \neq 0$.

Algorithm 2 Min-AM for strongly convex quadratic optimization

Input: $x_0 \in \mathbb{R}^d, \beta_k > 0, 0 < \max_iter \leq d$

```
1:  $p_0, q_0 = \mathbf{0} \in \mathbb{R}^d$ 
2: for  $k = 0, 1, \dots, \max\_iter$  do
3:    $r_k = -\nabla f(x_k)$ 
4:   if  $k > 0$  then
5:      $p = x_k - x_{k-1}, q = r_k - r_{k-1}$  (Compute  $\Delta x_{k-1}, \Delta r_{k-1}$ )
6:      $\zeta_k = (p_{k-1}^\top q_{k-1})^\dagger p_{k-1}^\top q$ 
7:      $p_k = p - p_{k-1}\zeta_k, q_k = q - q_{k-1}\zeta_k$  ( $p_k \perp q_{k-1}, q_k \perp p_{k-1}$ )
8:   end if
9:    $\Gamma_k^{(1)} = (p_k^\top q_k)^\dagger p_k^\top r_k$ 
10:   $x_k^{(1)} = x_k - p_k \Gamma_k^{(1)}, r_k^{(1)} = r_k - q_k \Gamma_k^{(1)}$  (Projection step:  $r_k^{(1)} \perp p_k$ )
11:   $x_k^{(2)} = x_k^{(1)} + \beta_k r_k^{(1)}$  (Mixing step)
12:   $\Gamma_k^{(2)} = \beta_k (p_k^\top q_k)^\dagger q_k^\top r_k^{(1)}$ 
13:   $x_{k+1} = x_k^{(2)} - p_k \Gamma_k^{(2)}$  (Projection step:  $r_{k+1} \perp p_k$ )
14:  if  $\|r_k^{(1)}\|_2 = 0$  then
15:    break
16:  end if
17: end for
18: return  $x_k$ 
```

We first prove that

$$q_k = -Ap_k \quad (42)$$

by induction. By the definition of f and $r_k = -\nabla f(x_k)$, it is clear that $r_k - r_{k-1} = -A(x_k - x_{k-1})$. So (42) holds for $k = 1$. For $k \geq 2$, suppose that (42) holds for $k-1$. Then $q_k = \Delta r_{k-1} - q_{k-1}\zeta_k = -A\Delta x_{k-1} + Ap_{k-1}\zeta_k = -Ap_k$, which completes the induction.

Suppose that $p_k^\top q_k \neq 0$. We have $\Gamma_k^{(1)} = (p_k^\top q_k)^{-1} p_k^\top r_k$. For $\Gamma_k^{(2)}$, the exact solution is

$$\begin{aligned} \Gamma_k^{(2)} &= (p_k^\top q_k)^{-1} p_k^\top r_k^{(2)} = (p_k^\top q_k)^{-1} p_k^\top (r_k^{(1)} - \beta_k A r_k^{(1)}) \\ &= -\beta_k (p_k^\top q_k)^{-1} p_k^\top A r_k^{(1)} = -\beta_k (p_k^\top q_k)^{-1} (A p_k)^\top r_k^{(1)} = \beta_k (p_k^\top q_k)^{-1} q_k^\top r_k^{(1)} \end{aligned} \quad (43)$$

due to $p_k^\top r_k^{(1)} = 0$, the symmetry of A , and (42). As a result,

$$\begin{aligned} x_{k+1} &= x_k^{(2)} - p_k \Gamma_k^{(2)} = x_k^{(1)} + \beta_k r_k^{(1)} - p_k \Gamma_k^{(2)} \\ &= x_k - p_k \Gamma_k^{(1)} + \beta_k r_k^{(1)} - p_k \Gamma_k^{(2)} \\ &= x_k - p_k \Gamma_k^{(1)} + \beta_k r_k^{(1)} - p_k \cdot \beta_k (p_k^\top q_k)^{-1} q_k^\top r_k^{(1)} \\ &= x_k - p_k (p_k^\top q_k)^{-1} p_k^\top r_k + \beta_k (I - p_k (p_k^\top q_k)^{-1} q_k^\top) r_k^{(1)} \\ &= x_k - p_k (p_k^\top q_k)^{-1} p_k^\top r_k + \beta_k (I - p_k (p_k^\top q_k)^{-1} q_k^\top) (I - q_k (p_k^\top q_k)^{-1} p_k^\top) r_k. \end{aligned} \quad (44)$$

The corresponding inverse Hessian approximation is

$$H_k = -p_k (p_k^\top q_k)^{-1} p_k^\top + \beta_k (I - p_k (p_k^\top q_k)^{-1} q_k^\top) (I - q_k (p_k^\top q_k)^{-1} p_k^\top). \quad (45)$$

The reason to use the projection condition of AM-I. Note that the Galerkin's projection condition of AM-I is used to determine $\Gamma_k^{(1)}$ and $\Gamma_k^{(2)}$. On the other side, if the Galerkin's condition of the original AM is used to determine $\Gamma_k^{(1)}$ and $\Gamma_k^{(2)}$, the scheme (41) leads to

$$x_{k+1} = x_k - p_k (q_k^\top q_k)^{-1} q_k^\top r_k + \beta_k (I + p_k (q_k^\top q_k)^{-1} q_k^\top A) (I - q_k (q_k^\top q_k)^{-1} q_k^\top) r_k, \quad (46)$$

which is not a practical method since Aq_k needs to be explicitly computed.

B.2 Proof of Theorem 1

Recall that the strongly convex quadratic optimization is formulated as

$$\min_{x \in \mathbb{R}^d} f(x) := \frac{1}{2} x^T A x - b^T x, \quad (47)$$

where $A \in \mathbb{R}^{d \times d}$ is SPD, $b \in \mathbb{R}^d$. Solving (47) is equivalent to solving the SPD linear system

$$A x = b. \quad (48)$$

We first state the relationship of the AM-I with the full orthogonalization method [52] (FOM) in the following propositions. Let x_k^{FOM} and $r_k^{\text{FOM}} := b - A x_k^{\text{FOM}}$ denote the k -th FOM iterate and residual, respectively. Define $e^j := (1, 1, \dots, 1)^T \in \mathbb{R}^j$ for $j \geq 1$. The main results of the full-memory AM-I are stated in Proposition 1 and Proposition 2. (Proposition 1 is a known result in [60].)

Proposition 1 (General linear system). *For solving a general linear system $Ax = b$ with the full-memory AM-I, suppose that $\beta_k > 0$ and the fixed-point map is $g(x) = (I - A)x + b$. If the initial point of AM-I is $x_0 = x_0^{\text{FOM}}$, and $X_j^T R_j$ is nonsingular for $j = 1, \dots, k$, then the intermediate iterate \bar{x}_k satisfies $\bar{x}_k = x_k^{\text{FOM}}$.*

We give the proof here, which is similar to [63, Proof of Proposition 1], but applies to the Type-I AM.

Proof. The definition of the fixed-point map suggests that the residual $r_k = g(x_k) - x_k = b - A x_k$. From (5), for $j = 1, \dots, k$, the nonsingularity of $X_j^T R_j$ ensures that each Γ_j is uniquely determined. Thus the updates of AM-I are well defined.

Since $X_k^T R_k$ is nonsingular, we have $\text{rank}(X_k) = \text{rank}(R_k) = k$. We first show

$$\text{range}(X_k) = \mathcal{K}_k(A, r_0^{\text{FOM}}) \quad (49)$$

by induction. We abbreviate $\mathcal{K}_k(A, r_0^{\text{FOM}})$ as \mathcal{K}_k in this proof.

First, $\Delta x_0 = \beta_0 r_0 = \beta_0 r_0^{\text{FOM}}$ since $x_1 = x_0 + \beta_0 r_0$. If $k = 1$, then the proof is complete. Then, suppose that $k > 1$ and, as an inductive hypothesis, that $\text{range}(X_{k-1}) = \mathcal{K}_{k-1}$. From (3), we have $x_{k+1} = x_k + \beta_k r_k - (X_k + \beta_k R_k) \Gamma_k$. Noting that $R_k = -A X_k$, it follows that

$$\begin{aligned} \Delta x_{k-1} &= x_k - x_{k-1} \\ &= \beta_{k-1} r_{k-1} - (X_{k-1} + \beta_{k-1} R_{k-1}) \Gamma_{k-1} \\ &= \beta_{k-1} (b - A x_{k-1}) - (X_{k-1} - \beta_{k-1} A X_{k-1}) \Gamma_{k-1} \\ &= \beta_{k-1} b - \beta_{k-1} A (x_0 + \Delta x_0 + \dots + \Delta x_{k-2}) - (X_{k-1} - \beta_{k-1} A X_{k-1}) \Gamma_{k-1} \\ &= \beta_{k-1} r_0 - \beta_{k-1} A X_{k-1} e^{k-1} - (X_{k-1} - \beta_{k-1} A X_{k-1}) \Gamma_{k-1}. \end{aligned} \quad (50)$$

Since $r_0 \in \mathcal{K}_{k-1}$, and by the inductive hypothesis $\text{range}(X_{k-1}) = \mathcal{K}_{k-1}$ which also implies $\text{range}(A X_{k-1}) \subseteq \mathcal{K}_k$, we know $\Delta x_{k-1} \in \mathcal{K}_k$. Thus, $\text{range}(X_k) \subseteq \mathcal{K}_k$. Since $\text{rank}(X_k) = k$, namely $\dim(\text{range}(X_k)) = \dim(\mathcal{K}_k)$, we have $\text{range}(X_k) = \mathcal{K}_k$, thus completing the induction. As a result, we also have

$$\text{range}(R_k) = \text{range}(A X_k) = A \mathcal{K}_k(A, r_0). \quad (51)$$

Recalling that to determine Γ_k , we solve the projection condition:

$$\bar{r}_k = r_k - R_k \Gamma_k \perp \text{range}(X_k). \quad (52)$$

The nonsingularity of $X_k^T R_k$ ensures that Γ_k is uniquely determined by

$$\Gamma_k = (X_k^T R_k)^{-1} X_k^T r_k. \quad (53)$$

Also, since $r_k = b - A x_k = b - A(x_0 + X_k e^k) = r_0 - A X_k e^k$, we have $r_k - R_k \Gamma_k = r_k + A X_k \Gamma_k = r_0 - A X_k e^k + A X_k \Gamma_k = r_0 - A X_k \tilde{\Gamma}$, where $\tilde{\Gamma} = e^k - \Gamma$, for $\forall \Gamma \in \mathbb{R}^k$. So Γ_k solves (52) if and only if $\tilde{\Gamma}_k = e^k - \Gamma_k$ solves

$$r_0 - A X_k \tilde{\Gamma}_k \perp \text{range}(X_k). \quad (54)$$

According to (49), the condition (54) is equivalent to

$$r_0 - Az \perp \text{range}(X_k) \text{ s.t. } z \in \mathcal{K}_k, \quad (55)$$

where $\text{range}(X_k) = \mathcal{K}_k$. Since the initializations are identical, the condition (55) for AM-I is the Petrov-Galerkin condition for FOM. Due to the nonsingularity of $X_k^T R_k$, the solution of (54) is also unique. Therefore, we have

$$\bar{x}_k = x_k - X_k \Gamma_k = x_k - X_k(e^k - \tilde{\Gamma}_k) = x_0 + X_k \tilde{\Gamma}_k = x_k^{\text{FOM}}.$$

□

In Proposition 1, the assumption that $X_k^T R_k$ is nonsingular is critical to ensure no stagnation occurs at the k -th iteration for solving a general linear system. In fact, for SPD linear systems (48) or strongly convex quadratic optimization (47), when AM-I breaks down, i.e. $X_k^T R_k$ is singular, AM-I obtains the exact solution, as shown in the next proposition.

Proposition 2 (SPD). *For applying the AM-I to minimize a strongly convex quadratic problem (47), or equivalently, solve an SPD linear system (48), suppose that $\beta_k > 0$ and the fixed-point map is $g(x) = (I - A)x + b$. If the condition $X_k^T R_k$ is nonsingular holds for $1 \leq k < s$ while failing to hold for $k = s$, where $s \geq 1$, then the residual of AM-I satisfies $r_s = \bar{r}_{s-1} = 0$.*

We give the proof here, which is similar to [63, Proof of Proposition 2], but applies to the Type-I AM.

Proof. The definition of g suggests that the residual $r_k = g(x_k) - x_k = b - Ax_k$. The relation $R_k = -AX_k$ holds during the iterations and the nonsingularity of A implies $\text{rank}(X_k) = \text{rank}(R_k)$. Since A is SPD and $X_k^T R_k = -X_k^T A X_k$, it follows that $X_k^T R_k$ being nonsingular $\Leftrightarrow \text{rank}(X_k) = \text{rank}(R_k) = k$. Hence $\text{rank}(X_k) = k$ holds for $1 \leq k < s$ while failing to hold for $k = s$.

For $s = 1$, since the first step of AM is $x_1 = x_0 + \beta_0 r_0$, the assumption $\text{rank}(X_1) = 0$ implies that $\text{rank}(r_0) = \text{rank}(X_1) = 0$, then $r_1 = \bar{r}_0 := 0$.

For $s > 1$, $\Delta x_{s-1} = x_s - x_{s-1} = -X_{s-1} \Gamma_{s-1} + \beta_{s-1} \bar{r}_{s-1}$. The rank deficiency of X_s implies $\Delta x_{s-1} \in \text{range}(X_{s-1})$, which further implies $\bar{r}_{s-1} \in \text{range}(X_{s-1})$. So there exists $\zeta \in \mathbb{R}^{s-1}$, such that $\bar{r}_{s-1} = X_{s-1} \zeta$. Due to $\bar{r}_{s-1} \perp X_{s-1}$, we have

$$0 = (\bar{r}_{s-1})^T X_{s-1} = (X_{s-1} \zeta)^T X_{s-1} = \zeta^T X_{s-1}^T X_{s-1}. \quad (56)$$

Because $\text{rank}(X_{s-1}) = s - 1$, we have $\zeta = 0$, which implies $\bar{r}_{s-1} = 0$. Hence $x_s = \bar{x}_{s-1}$ and $r_s = \bar{r}_{s-1} = 0$. □

It is also known that FOM method can be simplified to the conjugate gradient (CG) method [53] in this case. Now we prove Theorem 1, following a similar procedure of [63, Proof of Theorem 1].

Proof of Theorem 1. The A -norm minimization problem in the property (iii) is equivalent to the Galerkin condition [53], i.e. $z_k = \arg \min_{z \in \mathcal{K}_k(A, r_0)} \|x_0 + z - x^*\|_A \Leftrightarrow r_0 - Az_k \perp \mathcal{K}_k(A, r_0)$. Also, $r_k^{(3)} = r_{k+1} = r_k^{(2)} - q_k \Gamma_k^{(2)}$.

Besides relations (i)-(iii), we add an auxiliary relation here:

(iv) $r_k = r_0 + Q_k \bar{\Gamma}_k \in \mathcal{K}_{k+1}(A, r_0)$, where $\bar{\Gamma}_k \in \mathbb{R}^k$.

We prove the relations (i)-(iv) by induction.

For $k = 1$, since $r_0 = r_0^{(1)} \neq 0$, we have $\text{rank}(\Delta x_0) = \text{rank}(X_1) = 1$ and $\text{rank}(\Delta r_0) = \text{rank}(R_1) = 1$. So $p_k \neq 0$ and $p_k^T q_k = -p_k^T A p_k < 0$. The relation (i) holds. Since $p_1 = \Delta x_0$, $q_1 = \Delta r_0$, and $\Delta r_0 = -A \Delta x_0$, the equality $Q_1 = -A P_1$ also holds. Since $r_1 = r_0 - \beta_0 A r_0$ and $\text{range}(Q_1) = \text{span}\{A r_0\}$, it is clear that $r_1 = r_0 - Q_1 \bar{\Gamma}_1 \in \mathcal{K}_2(A, r_0)$, namely relation (iv). Due to the projection step (Line 10 in Algorithm 2), $r_1^{(1)} \perp \text{range}(P_1)$ and $r_1^{(1)}$ is unique, which is guaranteed by $p_1^T q_1 \neq 0$. Also, $r_1^{(1)} = r_1 - q_1 \Gamma_1^{(1)} = r_0 - \beta_0 A r_0 - q_1 \Gamma_1^{(1)} = r_0 - Q_1 \eta_1$, where the last equality is due to $\text{span}\{A r_0\} = \text{range}(Q_1)$. For $r_1^{\text{FOM}} = r_0 - A z_1$, where $z_1 \in \mathcal{K}_1(A, r_0)$, it holds $r_1^{\text{FOM}} \perp \mathcal{K}_1(A, r_0) = \text{range}(P_1)$. As a result, both $r_1^{(1)}$ and r_1^{FOM} are the oblique projections of r_0 onto the subspace $\text{range}(P_1)^\perp$ along $\text{range}(Q_1)$. Since $p_1^T q_1 \neq 0$, the projection exists and is

unique, which implies $r_1^{(1)} = r_1^{\text{FOM}}$. So $x_1^{(1)} = x_1^{\text{FOM}} = x_0 + z_1$ because their residuals are equal and A is nonsingular. Hence the relation (iii) holds.

Suppose that $k > 1$, and as an inductive hypothesis, the relations (i)-(iv) hold for $j = 1, \dots, k-1$. Consider the k -th iteration. From Line 7 in Algorithm 2, $q_k \in \text{span}\{\Delta r_{k-1}, q_{k-1}\}$, and $p_k \in \text{span}\{\Delta x_{k-1}, p_{k-1}\}$. Due to $p_{k-1}^T q_{k-1} = -p_{k-1}^T A p_{k-1} \neq 0$, we have $q_k^T p_{k-1} = 0$. We first prove that $p_k \neq 0$ by contradiction.

If $p_k = 0$, then from Line 7 in Algorithm 2, $\Delta x_{k-1} \in \text{span}\{p_{k-1}\}$. From Line 10, Line 11, and Line 13, we have

$$\Delta x_{k-1} = x_k - x_{k-1} = \beta_{k-1} r_{k-1}^{(1)} - p_{k-1} (\Gamma_{k-1}^{(1)} + \Gamma_{k-1}^{(2)}). \quad (57)$$

So $r_{k-1}^{(1)} \in \text{span}\{p_{k-1}\}$ since $\Delta x_{k-1} \in \text{span}\{p_{k-1}\}$. Hence there exists $\zeta \in \mathbb{R}$, such that $r_{k-1}^{(1)} = p_{k-1} \zeta$. From the Line 10, we know $r_{k-1}^{(1)} \perp \text{span}\{p_{k-1}\}$, so we have

$$0 = (r_{k-1}^{(1)})^T p_{k-1} = (p_{k-1} \zeta)^T p_{k-1} = \zeta^T p_{k-1}^T p_{k-1}.$$

Since $p_{k-1} \neq 0$, it follows that $\zeta = 0$ which implies $r_{k-1}^{(1)} = 0$. It is impossible otherwise Algorithm 2 has terminated in the $(k-1)$ -th iteration. So $p_k \neq 0$. Moreover $p_k^T q_k = -p_k^T A p_k \neq 0$ since A is SPD.

Since $r_{k-1}^{(1)} = r_{k-1} - q_{k-1} \Gamma_{k-1}^{(1)}$, and $r_{k-1} \in \mathcal{K}_k(A, r_0)$, $q_{k-1} \in \text{range}(Q_{k-1}) = A \mathcal{K}_{k-1}(A, r_0)$ due to the inductive hypothesis, we have $r_{k-1}^{(1)} \in \mathcal{K}_k(A, r_0)$, which together with (57) and $p_{k-1} \in \text{range}(P_{k-1}) = \mathcal{K}_{k-1}(A, r_0)$ infers $\Delta x_{k-1} \in \mathcal{K}_k(A, r_0)$. As a result, $p_k \in \text{span}\{\Delta x_{k-1}, p_{k-1}\} \subseteq \mathcal{K}_k(A, r_0)$. So $\text{range}(P_k) = \text{range}(P_{k-1}, p_k) \subseteq \mathcal{K}_k(A, r_0)$. Moreover, $p_k \notin \text{range}(P_{k-1})$. We prove it by contradiction.

If $p_k \in \text{range}(P_{k-1})$, then from Line 7 in Algorithm 2, $\Delta x_{k-1} \in \text{span}\{p_{k-1}, p_k\} \subseteq \text{range}(P_{k-1})$, which together with (57) leads to $r_{k-1}^{(1)} \in \text{range}(P_{k-1})$. Hence there exists $\xi \in \mathbb{R}^{k-1}$ such that $r_{k-1}^{(1)} = P_{k-1} \xi$. From the inductive hypothesis, we know $r_{k-1}^{(1)} \perp \text{range}(P_{k-1})$. So we have

$$0 = (r_{k-1}^{(1)})^T P_{k-1} = (P_{k-1} \xi)^T P_{k-1} = \xi^T P_{k-1}^T P_{k-1}.$$

Since $-P_{k-1}^T A P_{k-1} = P_{k-1}^T Q_{k-1}$ is nonsingular due to the inductive hypothesis $p_i \perp q_j$ ($1 \leq i \neq j \leq k-1$) and $p_i^T q_i = -p_i^T A p_i \neq 0$ ($1 \leq i \leq k-1$) as $p_i \neq 0$, it follows that $\text{rank}(P_{k-1}) = k-1$. So $P_{k-1}^T P_{k-1}$ is nonsingular which implies that $\xi = 0$. Then $r_{k-1}^{(1)} = P_{k-1} \xi = 0$. But it is impossible otherwise Algorithm 2 has terminated in the $(k-1)$ -th iteration. Hence $p_k \notin \text{range}(P_{k-1})$. As a result, $\text{range}(P_k) = \mathcal{K}_k(A, r_0)$.

Because $\Delta r_{k-1} = -A \Delta x_{k-1}$ and $q_{k-1} = -A p_{k-1}$, Line 7 in Algorithm 2 infers $q_k = -A p_k$. So $Q_k = -A P_k$. Hence $\text{range}(Q_k) = A \mathcal{K}_k(A, r_0)$.

To prove $p_i \perp q_j$ for $1 \leq i \neq j \leq k$, it suffices to show $q_k \perp P_{k-1}$. From the construction of q_k in Line 7 in Algorithm 2 and $p_{k-1}^T q_{k-1} \neq 0$, we know $q_k^T p_{k-1} = 0$. To further prove $q_k \perp \text{range}(P_{k-2})$ ($k \geq 3$), note that

$$\begin{aligned} \Delta r_{k-1} &= -A \Delta x_{k-1} = A p_{k-1} (\Gamma_{k-1}^{(1)} + \Gamma_{k-1}^{(2)}) - \beta_{k-1} A r_{k-1}^{(1)} \\ &= -q_{k-1} (\Gamma_{k-1}^{(1)} + \Gamma_{k-1}^{(2)}) - \beta_{k-1} A r_{k-1}^{(1)}, \end{aligned}$$

where the second equality is a direct substitution with (57). Therefore,

$$P_{k-2}^T \Delta r_{k-1} = -P_{k-2}^T q_{k-1} (\Gamma_{k-1}^{(1)} + \Gamma_{k-1}^{(2)}) - \beta_{k-1} P_{k-2}^T A r_{k-1}^{(1)} = 0 - \beta_{k-1} (A P_{k-2})^T r_{k-1}^{(1)} = 0, \quad (58)$$

where the second equality is due to $q_{k-1} \perp \text{range}(P_{k-2})$ and A is SPD, and the third equality is due to $r_{k-1}^{(1)} \perp \text{range}(P_{k-1}) = \mathcal{K}_{k-1}(A, r_0)$ and $\text{range}(A P_{k-2}) = A \mathcal{K}_{k-2}(A, r_0) \subseteq \mathcal{K}_{k-1}(A, r_0)$. As a result, we obtain

$$P_{k-2}^T q_k = P_{k-2}^T (\Delta r_{k-1} - q_{k-1} \zeta_k) = 0,$$

which is due to (58) and $q_{k-1} \perp \text{range}(P_{k-2})$. Therefore, $q_k \perp \text{range}(P_{k-1})$. Since A is SPD, we also have $p_k^T Q_{k-1} = -p_k^T A P_{k-1} = -(A p_k)^T P_{k-1} = q_k^T P_{k-1} = 0$. Hence relation (ii) holds in the k -th iteration.

Next, we prove the relation (iv). We have

$$\begin{aligned} r_k &= r_{k-1}^{(2)} - q_{k-1} \Gamma_{k-1}^{(2)} \\ &= r_{k-1}^{(1)} - \beta_{k-1} A r_{k-1}^{(1)} - q_{k-1} \Gamma_{k-1}^{(2)} \\ &= r_{k-1} - q_{k-1} \Gamma_{k-1}^{(1)} - \beta_{k-1} A (r_{k-1} - q_{k-1} \Gamma_{k-1}^{(1)}) - q_{k-1} \Gamma_{k-1}^{(2)} \\ &= r_0 + Q_{k-1} \bar{\Gamma}_{k-1} - q_{k-1} \Gamma_{k-1}^{(1)} - \beta_{k-1} A (r_0 + Q_{k-1} \bar{\Gamma}_{k-1}) + \beta_{k-1} A q_{k-1} \Gamma_{k-1}^{(1)} - q_{k-1} \Gamma_{k-1}^{(2)}, \end{aligned}$$

where the last equality is due to $r_{k-1} = r_0 + Q_{k-1} \bar{\Gamma}_{k-1}$ by the inductive hypothesis. Since $\text{range}(Q_{k-1}) = A \mathcal{K}_{k-1}(A, r_0) \subseteq A \mathcal{K}_k(A, r_0)$, $q_{k-1} \in \text{range}(Q_{k-1})$, $A r_0 \in A \mathcal{K}_k(A, r_0)$, $\text{range}(A Q_{k-1}) \subseteq A^2 \mathcal{K}_{k-1}(A, r_0) \subseteq A \mathcal{K}_k(A, r_0)$, and $A q_{k-1} \in \text{range}(A Q_{k-1})$, it is clear that $r_k = r_0 + Q_k \bar{\Gamma}_k \in \mathcal{K}_{k+1}(A, r_0)$ for some $\bar{\Gamma}_k \in \mathbb{R}^k$. The relation (iv) is proved.

Finally, we prove the relation (iii). For proving $r_k^{(1)} \perp \text{range}(P_k)$, note that $r_k^{(1)} \perp \text{span}\{p_k\}$ already holds due to the projection step (Line 9 and Line 10 in Algorithm 2) and $p_k^T q_k \neq 0$. It suffices to prove $r_k^{(1)} \perp \text{range}(P_{k-1})$. We first prove $r_k \perp \text{range}(P_{k-1})$. The projection step (Line 12 and Line 13 in Algorithm 2) at the $(k-1)$ -iteration ensures $r_k \perp p_{k-1}$. It remains to prove $r_k \perp \text{range}(P_{k-2})$. Because $r_k = r_{k-1}^{(2)} - q_{k-1} \Gamma_{k-1}^{(2)} = r_{k-1}^{(1)} - \beta_{k-1} A r_{k-1}^{(1)} - q_{k-1} \Gamma_{k-1}^{(2)}$, and $r_{k-1}^{(1)} \perp \text{range}(P_{k-2})$, $q_{k-1} \perp \text{range}(P_{k-2})$, $(A r_{k-1}^{(1)})^T P_{k-2} = (r_{k-1}^{(1)})^T (A P_{k-2}) = 0$ as $\text{range}(A P_{k-2}) \subseteq \mathcal{K}_{k-1}(A, r_0) = \text{range}(P_{k-1})$, it follows that $r_k \perp \text{range}(P_{k-2})$. Therefore, $r_k \perp \text{range}(P_{k-1})$. Also noting that $q_k \perp \text{range}(P_{k-1})$, we have $r_k^{(1)} = r_k - q_k \Gamma_k^{(1)} \perp \text{range}(P_{k-1})$. So $r_k^{(1)} \perp \text{range}(P_k)$.

To prove $x_k^{(1)} = x_k^{\text{FOM}} := x_0 + z_k$, where $z_k \in \mathcal{K}_k(A, r_0)$ should satisfy $r_0 - A z_k \perp \mathcal{K}_k(A, r_0)$, first we have $r_k = r_0 + Q_k \bar{\Gamma}_k$, where $\bar{\Gamma}_k \in \mathbb{R}^k$. Hence $r_k^{(1)} = r_k - q_k \Gamma_k^{(1)} = r_0 + Q_k \bar{\Gamma}_k - q_k \Gamma_k^{(1)} = r_0 - Q_k \eta_k$, where $\eta_k \in \mathbb{R}^k$. Since $P_k^T Q_k$ is nonsingular, there exists a unique $\eta_k \in \mathbb{R}^k$ such that $r_k^{(1)} \perp \text{range}(P_k)$, i.e., $r_k^{(1)}$ is the oblique projection of r_0 onto the subspace $\text{range}(P_k)^\perp$ along $\text{range}(Q_k)$. On the other side, for FOM, $r_k^{\text{FOM}} = r_0 - A z_k \perp \mathcal{K}_k(A, r_0) = \text{range}(P_k)$, so r_k^{FOM} is also the oblique projection of r_0 onto the subspace $\text{range}(P_k)^\perp$ along $\text{range}(Q_k)$. So $r_k^{(1)} = r_k^{\text{FOM}}$, which further indicates $x_k^{(1)} = x_k^{\text{FOM}}$. Hence, the relation (iii) holds.

With relations (i)-(iv) being proved in the k -th iteration, we complete the induction. \square

B.3 Convergence analysis

Since the basic Min-AM is equivalent to FOM (or CG) method in the sense that $\bar{x}_k^{(1)} = x_k^{\text{FOM}} = x_k^{\text{CG}}$ for strongly convex quadratic optimization, where x_k^{CG} is the k -th iterate of CG, we can obtain the convergence result as a corollary of Theorem 1.

Corollary 1. *For solving the strongly convex quadratic problem (47), let x^* be the exact solution, and $\{x_k\}$ is the sequence generated by the basic Min-AM described in Algorithm 2. Define $\theta_k = \|I - \beta_k A\|_2$. Then*

$$\|x_{k+1} - x^*\|_A \leq \theta_k \min_{p \in \mathcal{P}_k, p(0)=1} \|p(A)(x_0 - x^*)\|_A, \quad (59)$$

where \mathcal{P}_k denotes the space of polynomials of degree not exceeding k . Moreover, the algorithm finds the exact solution in at most $(d+1)$ iterations.

Proof. From Theorem 1, the classical convergence analysis of the CG method [28, 53], and assuming that CG and Min-AM are initialized from the same starting point, i.e., $x_0^{\text{CG}} = x_0$, we have

$$\|x_k^{(1)} - x^*\|_A = \|x_k^{\text{CG}} - x^*\|_A = \min_{p \in \mathcal{P}_k, p(0)=1} \|p(A)(x_0 - x^*)\|_A. \quad (60)$$

Since $x_k^{(2)} = x_k^{(1)} + \beta_k r_k^{(1)}$, it follows that

$$x_k^{(2)} - x^* = x_k^{(1)} - x^* + \beta_k(b - Ax_k^{(1)}) = x_k^{(1)} - x^* - \beta_k A(x_k^{(1)} - x^*) = (I - \beta_k A)(x_k^{(1)} - x^*).$$

Hence,

$$\|x_k^{(2)} - x^*\|_A \leq \|I - \beta_k A\|_A \|x_k^{(1)} - x^*\|_A = \theta_k \|x_k^{(1)} - x^*\|_A, \quad (61)$$

where we use the fact that

$$\|I - \beta_k A\|_A = \max_{\|x\|_A=1} \|(I - \beta_k A)x\|_A = \max_{\|A^{1/2}x\|_2=1} \|A^{1/2}(I - \beta_k A)x\|_2 \quad (62)$$

$$= \max_{\|y\|_2=1} \|A^{1/2}(I - \beta_k A)A^{-1/2}y\|_2 = \|A^{1/2}(I - \beta_k A)A^{-1/2}\|_2 = \|I - \beta_k A\|_2. \quad (63)$$

Noting that $x_{k+1} = x_k^{(2)} - p_k \Gamma_k^{(2)}$, $r_{k+1} = r_k^{(2)} + Ap_k \Gamma_k^{(2)} \perp p_k$, it follows that

$$\Gamma_k^{(2)} = \arg \min_{\Gamma \in \mathbb{R}} \|x_k^{(2)} - p_k \Gamma - x^*\|_A.$$

So

$$\|x_{k+1} - x^*\|_A \leq \|x_k^{(2)} - x^*\|_A. \quad (64)$$

Combining (60), (61), and (64) yields (59).

Since CG iterates at most d steps to obtain x^* , the Min-AM will obtain the same solution in at most $(d + 1)$ steps. \square

C Details of the restarted Min-AM

C.1 Proof of Theorem 2

For the proof of Theorem 2 about the restarted Min-AM, we use a local quadratic model inspired by the techniques in [17, 18]. The second-order Taylor expansion \hat{f} of f around x^* is

$$\hat{f}(x) = f(x^*) + \frac{1}{2}(x - x^*)^T A(x - x^*), \quad (65)$$

where $A := \nabla^2 f(x^*)$.

We compare the iterates $\{x_k\}$ generated by Min-AM to the iterates $\{\hat{x}_k\}$ generated by applying Min-AM to minimizing the quadratic model (65). More precisely, we give the following definition.

Definition 1. Let the mixing parameter $\{\beta_k\}$ be chosen to satisfy $\beta \leq |\beta_k| \leq \beta'$, where β and β' are two positive constants. The sequences $\{x_k\}$ and $\{\hat{x}_k\}$ are generated by the following two processes:

(1) *Process I:* Solve the optimization problem (1) with the restarted Min-AM (see Algorithm 1), and the resulting sequence is $\{x_k\}$.

(2) *Process II:* Apply the basic Min-AM to minimize $\hat{f}(x)$ in each interval between two successive restarts in Process I. Specifically, let m_k, β_k be the same as those in Process I. Define $\hat{r}_k = -\nabla \hat{f}(\hat{x}_k)$. Then, $\hat{x}_k = x_k$, and $\hat{p}_k = \hat{q}_k = \mathbf{0} \in \mathbb{R}^d$, if $m_k = 0$. For $m_k \geq 1$,

$$\Delta \hat{x}_{k-1} = \hat{x}_k - \hat{x}_{k-1}, \quad \Delta \hat{r}_{k-1} = \hat{r}_k - \hat{r}_{k-1},$$

$$\hat{p}_k = \Delta \hat{x}_{k-1} - \hat{p}_{k-1} \hat{\zeta}_k, \quad \hat{q}_k = \Delta \hat{r}_{k-1} - \hat{q}_{k-1} \hat{\zeta}_k, \quad \text{where } \hat{\zeta}_k = (\hat{p}_{k-1}^T \hat{q}_{k-1})^\dagger \hat{p}_{k-1}^T \Delta \hat{r}_{k-1};$$

defining $\hat{r}_k^{(1)} = \hat{r}_k - \hat{q}_k \hat{\Gamma}_k^{(1)}$, the update of \hat{x}_k is

$$\begin{aligned} \hat{x}_k^{(1)} &= \hat{x}_k - \hat{p}_k \hat{\Gamma}_k^{(1)}, \quad \text{where } \hat{\Gamma}_k^{(1)} = (\hat{p}_k^T \hat{q}_k)^\dagger \hat{p}_k^T \hat{r}_k, \\ \hat{x}_k^{(2)} &= \hat{x}_k^{(1)} + \beta_k \hat{r}_k^{(1)}, \\ \hat{x}_{k+1} &= \hat{x}_k^{(2)} - \hat{p}_k \hat{\Gamma}_k^{(2)}, \quad \text{where } \hat{\Gamma}_k^{(2)} = \beta_k (\hat{p}_k^T \hat{q}_k)^\dagger \hat{q}_k^T \hat{r}_k^{(1)}. \end{aligned} \quad (66)$$

By comparing $\{x_k\}$ and $\{\hat{x}_k\}$, we first have the following lemma.

Lemma 1. Suppose that Assumption 1 holds for the optimization problem (1). For the sequences $\{x_k\}$ and $\{\hat{x}_k\}$ defined in Definition 1, if $\|x_0 - x^*\|_2$ is sufficiently small and there exists a positive constant η_0 such that $\|\nabla f(x_j)\|_2 \leq \eta_0 \|\nabla f(x_0)\|_2$ for $j = 0, \dots, k$, then

$$\|r_k - \hat{r}_k\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2), \quad (67)$$

$$\|x_{k+1} - \hat{x}_{k+1}\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2). \quad (68)$$

Proof. Besides (67) and (68), we will also prove the following relations.

$$x_k \in \mathcal{B}_{\hat{\rho}}(x^*), \quad (69)$$

$$|\zeta_k| = \mathcal{O}(1), \quad |\hat{\zeta}_k - \zeta_k| = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2), \quad (70)$$

$$\|p_k\|_2 = \mathcal{O}(\|x_{k-m_k} - x^*\|_2), \quad \|q_k\|_2 = \mathcal{O}(\|x_{k-m_k} - x^*\|_2) \quad (71)$$

$$\|p_k - \hat{p}_k\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2), \quad \|q_k - \hat{q}_k\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2), \quad (72)$$

$$|\Gamma_k^{(1)}| = \mathcal{O}(1), \quad |\hat{\Gamma}_k^{(1)} - \Gamma_k^{(1)}| = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2), \quad (73)$$

$$\|x_k^{(1)} - \hat{x}_k^{(1)}\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2), \quad \|r_k^{(1)} - \hat{r}_k^{(1)}\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2), \quad (74)$$

$$\|x_k^{(2)} - \hat{x}_k^{(2)}\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2), \quad (75)$$

$$|\Gamma_k^{(2)}| = \mathcal{O}(1), \quad |\hat{\Gamma}_k^{(2)} - \Gamma_k^{(2)}| = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2). \quad (76)$$

Here, for convenience, we define $\zeta_k = \hat{\zeta}_k = 0$ if $m_k = 0$.

We first prove (69). Recall that $r_k = -\nabla f(x_k)$. Due to (12a), we have

$$\mu \|x_k - x^*\|_2 \leq \|r_k\|_2 = \|\nabla f(x_k) - \nabla f(x^*)\|_2 \leq L \|x_k - x^*\|_2. \quad (77)$$

Note that $\|r_k\|_2 \leq \eta_0 \|r_0\|_2$. By choosing $\|x_0 - x^*\|_2 \leq \frac{\mu \hat{\rho}}{\eta_0 L}$, we can ensure

$$\|x_k - x^*\|_2 \leq \frac{1}{\mu} \|r_k\|_2 \leq \frac{\eta_0}{\mu} \|r_0\|_2 \leq \frac{\eta_0 L}{\mu} \|x_0 - x^*\|_2 \leq \frac{\eta_0 L}{\mu} \cdot \frac{\mu \hat{\rho}}{\eta_0 L} = \hat{\rho}. \quad (78)$$

Then (69) holds. Since $\|r_j\|_2 \leq \eta_0 \|r_0\|_2$ holds for $j = 0, \dots, k$, (78) implies that a sufficiently small $\|x_0 - x^*\|_2$ ensures $\|x_{k-m_k} - x^*\|_2 \leq \frac{\eta_0 L}{\mu} \|x_0 - x^*\|_2$ is sufficiently small. Next, we prove (67), (68), and (70)-(76) by induction.

For $k = 0$, we have $\Gamma_k^{(1)} = \hat{\Gamma}_k^{(1)} = 0$, $\Gamma_k^{(2)} = \hat{\Gamma}_k^{(2)} = 0$. The relations (70)-(74), and (76) clearly hold. Also, due to (12b),

$$\|r_0 - \hat{r}_0\|_2 = \|\nabla f(x_0) - \nabla \hat{f}(\hat{x}_0)\|_2 = \|\nabla f(x_0) - \nabla^2 f(x^*)(x_0 - x^*)\|_2 \leq \frac{1}{2} \hat{\kappa} \|x_0 - x^*\|_2^2,$$

namely (67). Note that $x_{k+1}^{(2)} = x_k + \beta_k r_k$ and $\hat{x}_{k+1}^{(2)} = \hat{x}_k + \beta_k \hat{r}_k$ in this case. Then (68) and (75) follow from

$$\|x_1 - \hat{x}_1\|_2 = \|x_0 + \beta_0 r_0 - (\hat{x}_0 + \beta_0 \hat{r}_0)\|_2 = \beta_0 \|r_0 - \hat{r}_0\|_2 \leq \frac{\beta_0 \hat{\kappa}}{2} \|x_0 - x^*\|_2^2.$$

Suppose that $k \geq 1$, and as an inductive hypothesis, the relations (67), (68), and (70)-(76) hold for $j = 0, \dots, k-1$. Consider the k -th iteration.

If $m_k = 0$, i.e., a condition in (11a)-(11c) is violated at the beginning of the k -th iteration, then $\hat{x}_k = x_k$. The same as the case that $k = 0$, (67), (68), and (70)-(76) hold.

Consider the nontrivial case that $m_k \geq 1$. From (11c), it follows that

$$\|x_j - x^*\|_2 \leq \frac{1}{\mu} \|r_j\|_2 \leq \frac{\eta}{\mu} \|r_{k-m_k}\|_2 \leq \frac{\eta L}{\mu} \|x_{k-m_k} - x^*\|_2, \quad j = k - m_k + 1, \dots, k.$$

Therefore,

$$\|x_j - x^*\|_2 = \mathcal{O}(\|x_{k-m_k} - x^*\|_2), \quad j = k - m_k, \dots, k. \quad (79)$$

Since $x_k \in \mathcal{B}_{\hat{\rho}}(x^*)$, we have

$$\begin{aligned}
\|r_k - \hat{r}_k\|_2 &= \|\nabla f(x_k) - \nabla \hat{f}(\hat{x}_k)\|_2 \\
&\leq \|\nabla f(x_k) - \nabla \hat{f}(x_k)\|_2 + \|\nabla \hat{f}(x_k) - \nabla \hat{f}(\hat{x}_k)\|_2 \\
&= \|\nabla f(x_k) - \nabla^2 f(x^*)(x_k - x^*)\|_2 + \|\nabla^2 f(x^*)(x_k - \hat{x}_k)\|_2 \\
&\leq \frac{1}{2} \hat{\kappa} \|x_k - x^*\|_2^2 + L \|x_k - \hat{x}_k\|_2 \\
&= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2),
\end{aligned}$$

where the second inequality is due to (12b) and (12a), and the last inequality is due to (79) and the inductive hypothesis (68). Hence, the relation (67) holds. Next, we prove the relation (70).

Due to the check (11b), with $\underline{\kappa} := \tau \mu^3 \beta^2$, we have

$$\begin{aligned}
|p_j^T q_j| &\geq \tau |p_{k-m_k+1}^T q_{k-m_k+1}| = \tau |\Delta x_{k-m_k}^T \Delta r_{k-m_k}| \\
&\geq \tau \left| \Delta x_{k-m_k}^T \int_0^1 \nabla^2 f(x_{k-m_k} + t \Delta x_{k-m_k}) \Delta x_{k-m_k} dt \right| \\
&\geq \tau \mu \|\Delta x_{k-m_k}\|_2^2 = \tau \mu \beta_{k-m_k}^2 \|r_{k-m_k}\|_2^2 \\
&\geq \tau \mu^3 \beta^2 \|x_{k-m_k} - x^*\|_2^2 = \underline{\kappa} \|x_{k-m_k} - x^*\|_2^2,
\end{aligned} \tag{80}$$

for $j = k - m_k + 1, \dots, k$.

If $m_k = 1$, then $\zeta_k = \hat{\zeta}_k = 0$, and (70) holds. For $m_k > 1$, we have

$$|\zeta_k| = \left| \frac{p_{k-1}^T \Delta r_{k-1}}{p_{k-1}^T q_{k-1}} \right| \leq \frac{\|p_{k-1}\|_2 \cdot \|\Delta r_{k-1}\|_2}{\underline{\kappa} \|x_{k-m_k} - x^*\|_2^2} = \frac{\mathcal{O}(\|x_{k-m_k} - x^*\|_2^2)}{\underline{\kappa} \|x_{k-m_k} - x^*\|_2^2} = \mathcal{O}(1), \tag{81}$$

due to (71) and $\|\Delta r_{k-1}\|_2 \leq \|r_k\|_2 + \|r_{k-1}\|_2 \leq 2\eta \|r_{k-m_k}\|_2 = \mathcal{O}(\|x_{k-m_k} - x^*\|_2)$.

Next, if $p_{k-1}^T \Delta r_{k-1} \neq 0$, then

$$\begin{aligned}
|\zeta_k - \hat{\zeta}_k| &= |\zeta_k| \cdot \left| 1 - \frac{\hat{\zeta}_k}{\zeta_k} \right| = |\zeta_k| \cdot \left| 1 - \frac{\hat{p}_{k-1}^T \Delta \hat{r}_{k-1}}{\hat{p}_{k-1}^T \hat{q}_{k-1}} \cdot \frac{p_{k-1}^T q_{k-1}}{p_{k-1}^T \Delta r_{k-1}} \right| \\
&= |\zeta_k| \cdot \left| 1 - \frac{\hat{p}_{k-1}^T \Delta \hat{r}_{k-1}}{p_{k-1}^T \Delta r_{k-1}} \cdot \frac{p_{k-1}^T q_{k-1}}{\hat{p}_{k-1}^T \hat{q}_{k-1}} \right| \\
&= |\zeta_k| \cdot |a(1-b) + b| \leq |\zeta_k| \cdot (|a| + |b| + |ab|),
\end{aligned} \tag{82}$$

where $a := 1 - \frac{\hat{p}_{k-1}^T \Delta \hat{r}_{k-1}}{p_{k-1}^T \Delta r_{k-1}}$, $b := 1 - \frac{p_{k-1}^T q_{k-1}}{\hat{p}_{k-1}^T \hat{q}_{k-1}}$. We have

$$|\zeta_k| \cdot |a| = \left| \frac{p_{k-1}^T \Delta r_{k-1}}{p_{k-1}^T q_{k-1}} \right| \cdot \left| \frac{p_{k-1}^T \Delta r_{k-1} - \hat{p}_{k-1}^T \Delta \hat{r}_{k-1}}{p_{k-1}^T \Delta r_{k-1}} \right| = \left| \frac{p_{k-1}^T \Delta r_{k-1} - \hat{p}_{k-1}^T \Delta \hat{r}_{k-1}}{p_{k-1}^T q_{k-1}} \right|. \tag{83}$$

From (67) and (71),

$$|p_{k-1}^T (\Delta r_{k-1} - \Delta \hat{r}_{k-1})| \leq \|p_{k-1}\|_2 \cdot \|r_k - r_{k-1} - \hat{r}_k + \hat{r}_{k-1}\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3).$$

We also have

$$\begin{aligned}
|(p_{k-1} - \hat{p}_{k-1})^T \Delta \hat{r}_{k-1}| &= |(p_{k-1} - \hat{p}_{k-1})^T A \Delta \hat{x}_{k-1}| \\
&= |(p_{k-1} - \hat{p}_{k-1})^T A \Delta x_{k-1} - (p_{k-1} - \hat{p}_{k-1})^T A (\Delta x_{k-1} - \Delta \hat{x}_{k-1})| \\
&\leq |(p_{k-1} - \hat{p}_{k-1})^T A \Delta x_{k-1}| + |(p_{k-1} - \hat{p}_{k-1})^T A (\Delta x_{k-1} - \Delta \hat{x}_{k-1})| \\
&= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3) + \hat{\kappa}^2 \mathcal{O}(\|x_{k-m_k} - x^*\|_2^4) \\
&= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3).
\end{aligned}$$

Hence,

$$\begin{aligned}
|p_{k-1}^T \Delta r_{k-1} - \hat{p}_{k-1}^T \Delta \hat{r}_{k-1}| &= |p_{k-1}^T (\Delta r_{k-1} - \Delta \hat{r}_{k-1}) + (p_{k-1} - \hat{p}_{k-1})^T \Delta \hat{r}_{k-1}| \\
&\leq |p_{k-1}^T (\Delta r_{k-1} - \Delta \hat{r}_{k-1})| + |(p_{k-1} - \hat{p}_{k-1})^T \Delta \hat{r}_{k-1}| \\
&= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3).
\end{aligned} \tag{84}$$

Combining (83) and (84), we have

$$|\zeta_k| \cdot |a| \leq \frac{\hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3)}{\underline{\kappa} \|x_{k-m_k} - x^*\|_2^2} = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2). \tag{85}$$

Note that

$$\begin{aligned}
|\hat{p}_{k-1}^T \hat{q}_{k-1} - p_{k-1}^T q_{k-1}| &= |\hat{p}_{k-1}^T (\hat{q}_{k-1} - q_{k-1}) + (\hat{p}_{k-1} - p_{k-1})^T q_{k-1}| \\
&= |(\hat{p}_{k-1} - p_{k-1})^T (\hat{q}_{k-1} - q_{k-1}) + p_{k-1}^T (\hat{q}_{k-1} - q_{k-1}) + (\hat{p}_{k-1} - p_{k-1})^T q_{k-1}| \\
&\leq |(\hat{p}_{k-1} - p_{k-1})^T (\hat{q}_{k-1} - q_{k-1})| + |p_{k-1}^T (\hat{q}_{k-1} - q_{k-1})| + |(\hat{p}_{k-1} - p_{k-1})^T q_{k-1}| \\
&= \hat{\kappa}^2 \mathcal{O}(\|x_{k-m_k} - x^*\|_2^4) + \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3) + \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3) \\
&= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3),
\end{aligned} \tag{86}$$

and

$$\begin{aligned}
|\hat{p}_{k-1}^T \hat{q}_{k-1}| &\geq |p_{k-1}^T q_{k-1}| - |\hat{p}_{k-1}^T \hat{q}_{k-1} - p_{k-1}^T q_{k-1}| \\
&\geq \underline{\kappa} \|x_{k-m_k} - x^*\|_2^2 - \hat{\kappa} c_1 \|x_{k-m_k} - x^*\|_2^3 \geq \frac{1}{2} \underline{\kappa} \|x_{k-m_k} - x^*\|_2^2,
\end{aligned} \tag{87}$$

where the existence of the constant c_1 is guaranteed by (86), and the last inequality holds by choosing $\|x_0 - x^*\|_2 \leq \frac{\mu \underline{\kappa}}{2\eta_0 L \hat{\kappa} c_1}$, which ensures $\hat{\kappa} c_1 \|x_{k-m_k} - x^*\|_2 \leq \hat{\kappa} c_1 \frac{\eta_0 L}{\mu} \|x_0 - x^*\|_2 \leq \frac{1}{2} \underline{\kappa}$.

Then, we have

$$|b| = \left| 1 - \frac{p_{k-1}^T q_{k-1}}{\hat{p}_{k-1}^T \hat{q}_{k-1}} \right| = \left| \frac{\hat{p}_{k-1}^T \hat{q}_{k-1} - p_{k-1}^T q_{k-1}}{\hat{p}_{k-1}^T \hat{q}_{k-1}} \right| = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2). \tag{88}$$

As a result, by (85), (88), (81), and (82), we have

$$|\zeta_k - \hat{\zeta}_k| \leq |\zeta_k| |a| + |\zeta_k| |b| + |\zeta_k| |a| |b| = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2). \tag{89}$$

Now consider the case that $p_{k-1}^T \Delta r_{k-1} = 0$. It is clear that $\zeta_k = 0$. Then

$$|\zeta_k - \hat{\zeta}_k| = \left| \frac{\hat{p}_{k-1}^T \Delta \hat{r}_{k-1}}{\hat{p}_{k-1}^T \hat{q}_{k-1}} \right| \leq \frac{\hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3)}{\frac{1}{2} \underline{\kappa} \|x_{k-m_k} - x^*\|_2^2} = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2),$$

where the inequality is due to (84) and (87).

Hence the (70) holds.

For the relation (71), it holds since

$$\begin{aligned}
\|p_k\|_2 &= \|\Delta x_{k-1} - p_{k-1} \zeta_k\|_2 \leq \|x_k - x^*\|_2 + \|x_{k-1} - x^*\|_2 + \|p_{k-1}\|_2 |\zeta_k| \\
&= \mathcal{O}(\|x_{k-m_k} - x^*\|_2), \\
\|q_k\|_2 &= \|\Delta r_{k-1} - q_{k-1} \zeta_k\|_2 \leq L \|x_k - x_{k-1}\|_2 + \|q_{k-1}\|_2 |\zeta_k| \\
&\leq L \|x_k - x^*\|_2 + L \|x_{k-1} - x^*\|_2 + \|q_{k-1}\|_2 |\zeta_k| \\
&= \mathcal{O}(\|x_{k-m_k} - x^*\|_2).
\end{aligned}$$

Next, we prove (72). We have

$$\begin{aligned}
\|p_{k-1} \zeta_k - \hat{p}_{k-1} \hat{\zeta}_k\|_2 &= \|p_{k-1} \zeta_k - \hat{p}_{k-1} (\hat{\zeta}_k - \zeta_k) - \hat{p}_{k-1} \zeta_k\|_2 \\
&\leq \|(p_{k-1} - \hat{p}_{k-1}) \zeta_k\|_2 + \|\hat{p}_{k-1} (\hat{\zeta}_k - \zeta_k)\|_2 \\
&\leq \|(p_{k-1} - \hat{p}_{k-1}) \zeta_k\|_2 + \|(\hat{p}_{k-1} - p_{k-1}) (\hat{\zeta}_k - \zeta_k)\|_2 + \|p_{k-1} (\hat{\zeta}_k - \zeta_k)\|_2 \\
&= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2),
\end{aligned} \tag{90}$$

$$\begin{aligned}
& \|q_{k-1}\zeta_k - \hat{q}_{k-1}\hat{\zeta}_k\|_2 = \|q_{k-1}\zeta_k - \hat{q}_{k-1}(\hat{\zeta}_k - \zeta_k) - \hat{q}_{k-1}\zeta_k\|_2 \\
& \leq \|(q_{k-1} - \hat{q}_{k-1})\zeta_k\|_2 + \|(\hat{q}_{k-1} - q_{k-1})(\hat{\zeta}_k - \zeta_k)\|_2 + \|q_{k-1}(\hat{\zeta}_k - \zeta_k)\|_2 \\
& = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^2).
\end{aligned} \tag{91}$$

Therefore,

$$\begin{aligned}
\|p_k - \hat{p}_k\|_2 & = \|\Delta x_{k-1} - p_{k-1}\zeta_k - (\Delta \hat{x}_{k-1} - \hat{p}_{k-1}\hat{\zeta}_k)\|_2 \\
& \leq \|\Delta x_{k-1} - \Delta \hat{x}_{k-1}\|_2 + \|p_{k-1}\zeta_k - \hat{p}_{k-1}\hat{\zeta}_k\|_2 = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^2),
\end{aligned} \tag{92}$$

$$\begin{aligned}
\|q_k - \hat{q}_k\|_2 & = \|\Delta r_{k-1} - q_{k-1}\zeta_k - (\Delta \hat{r}_{k-1} - \hat{q}_{k-1}\hat{\zeta}_k)\|_2 \\
& \leq \|\Delta r_{k-1} - \Delta \hat{r}_{k-1}\|_2 + \|q_{k-1}\zeta_k - \hat{q}_{k-1}\hat{\zeta}_k\|_2 = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^2).
\end{aligned} \tag{93}$$

The (72) is proved.

Now, we prove (73), which follows a very similar procedure in proving (70). For the completeness, we still give the proof.

We have

$$|\Gamma_k^{(1)}| = \left| \frac{p_k^T r_k}{p_k^T q_k} \right| \leq \frac{\|p_k\|_2 \cdot \|r_k\|_2}{\underline{\kappa}\|x_{k-m_k} - x^*\|_2^2} = \frac{\mathcal{O}(\|x_{k-m_k} - x^*\|_2^2)}{\underline{\kappa}\|x_{k-m_k} - x^*\|_2^2} = \mathcal{O}(1), \tag{94}$$

due to (71) and $\|r_k\|_2 \leq \eta\|r_{k-m_k}\|_2 = \mathcal{O}(\|x_{k-m_k} - x^*\|_2)$.

Next, if $p_k^T r_k \neq 0$, then

$$\begin{aligned}
|\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)}| & = |\Gamma_k^{(1)}| \cdot \left| 1 - \frac{\hat{\Gamma}_k^{(1)}}{\Gamma_k^{(1)}} \right| = |\Gamma_k^{(1)}| \cdot \left| 1 - \frac{\hat{p}_k^T \hat{r}_k}{\hat{p}_k^T \hat{q}_k} \cdot \frac{p_k^T q_k}{p_k^T r_k} \right| = |\Gamma_k^{(1)}| \cdot \left| 1 - \frac{\hat{p}_k^T \hat{r}_k}{p_k^T r_k} \cdot \frac{p_k^T q_k}{\hat{p}_k^T \hat{q}_k} \right| \\
& = |\Gamma_k^{(1)}| \cdot |a_1(1 - b_1) + b_1| \leq |\Gamma_k^{(1)}| \cdot (|a_1| + |b_1| + |a_1 b_1|),
\end{aligned} \tag{95}$$

where $a_1 := 1 - \frac{\hat{p}_k^T \hat{r}_k}{p_k^T r_k}$, $b_1 := 1 - \frac{p_k^T q_k}{\hat{p}_k^T \hat{q}_k}$. We have

$$|\Gamma_k^{(1)}| \cdot |a_1| = \left| \frac{p_k^T r_k}{p_k^T q_k} \right| \cdot \left| \frac{p_k^T r_k - \hat{p}_k^T \hat{r}_k}{p_k^T r_k} \right| = \left| \frac{p_k^T r_k - \hat{p}_k^T \hat{r}_k}{p_k^T q_k} \right|. \tag{96}$$

From (67) and (71) which have been proved,

$$|p_k^T(r_k - \hat{r}_k)| \leq \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^3).$$

We also have

$$\begin{aligned}
|(p_k - \hat{p}_k)^T \hat{r}_k| & = |(p_k - \hat{p}_k)^T A(\hat{x}_k - x^*)| \\
& = |(p_k - \hat{p}_k)^T A(x_k - x^*) - (p_k - \hat{p}_k)^T A(x_k - \hat{x}_k)| \\
& \leq |(p_k - \hat{p}_k)^T A(x_k - x^*)| + |(p_k - \hat{p}_k)^T A(x_k - \hat{x}_k)| \\
& = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^3) + \hat{\kappa}^2\mathcal{O}(\|x_{k-m_k} - x^*\|_2^4) = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^3).
\end{aligned}$$

Hence,

$$\begin{aligned}
|p_k^T r_k - \hat{p}_k^T \hat{r}_k| & = |p_k^T(r_k - \hat{r}_k) + (p_k - \hat{p}_k)^T \hat{r}_k| \\
& \leq |p_k^T(r_k - \hat{r}_k)| + |(p_k - \hat{p}_k)^T \hat{r}_k| = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^3).
\end{aligned} \tag{97}$$

Combining (96) and (97), we have

$$|\Gamma_k^{(1)}| \cdot |a_1| \leq \frac{\hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^3)}{\underline{\kappa}\|x_{k-m_k} - x^*\|_2^2} = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2). \tag{98}$$

Note that

$$\begin{aligned}
|\hat{p}_k^T \hat{q}_k - p_k^T q_k| & = |\hat{p}_k^T(\hat{q}_k - q_k) + (\hat{p}_k - p_k)^T q_k| \\
& = |(\hat{p}_k - p_k)^T(\hat{q}_k - q_k) + p_k^T(\hat{q}_k - q_k) + (\hat{p}_k - p_k)^T q_k| \\
& \leq |(\hat{p}_k - p_k)^T(\hat{q}_k - q_k)| + |p_k^T(\hat{q}_k - q_k)| + |(\hat{p}_k - p_k)^T q_k| \\
& = \hat{\kappa}^2\mathcal{O}(\|x_{k-m_k} - x^*\|_2^4) + \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^3) + \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^3) \\
& = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^3),
\end{aligned} \tag{99}$$

and

$$\begin{aligned} |\hat{p}_k^T \hat{q}_k| &\geq |p_k^T q_k| - |\hat{p}_k^T \hat{q}_k - p_k^T q_k| \\ &\geq \underline{\kappa} \|x_{k-m_k} - x^*\|_2^2 - \hat{\kappa} c_2 \|x_{k-m_k} - x^*\|_2^3 \geq \frac{1}{2} \underline{\kappa} \|x_{k-m_k} - x^*\|_2^2, \end{aligned} \quad (100)$$

where the existence of the constant c_2 is guaranteed by (99), and the last inequality holds by choosing $\|x_0 - x^*\|_2 \leq \frac{\mu \underline{\kappa}}{2\eta_0 L \hat{\kappa} c_2}$ which ensures $\hat{\kappa} c_2 \|x_{k-m_k} - x^*\|_2 \leq \hat{\kappa} c_2 \frac{\eta_0 L}{\mu} \|x_0 - x^*\|_2 \leq \frac{1}{2} \underline{\kappa}$.

Then, we have

$$|b_1| = \left| 1 - \frac{p_k^T q_k}{\hat{p}_k^T \hat{q}_k} \right| = \left| \frac{\hat{p}_k^T \hat{q}_k - p_k^T q_k}{\hat{p}_k^T \hat{q}_k} \right| = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2). \quad (101)$$

As a result, by (98), (101), (94), and (95), we have

$$|\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)}| \leq |\Gamma_k^{(1)}| |a_1| + |\Gamma_k^{(1)}| |b_1| + |\Gamma_k^{(1)}| |a_1| |b_1| = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2). \quad (102)$$

Now consider the case that $p_k^T r_k = 0$. It is clear that $\Gamma_k^{(1)} = 0$. Then

$$|\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)}| = \left| \frac{\hat{p}_k^T \hat{r}_k}{\hat{p}_k^T \hat{q}_k} \right| \leq \frac{\hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3)}{\frac{1}{2} \underline{\kappa} \|x_{k-m_k} - x^*\|_2^2} = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2),$$

where the inequality is due to (97) and (100).

Hence the (73) holds.

Now, we prove (74). We have

$$\begin{aligned} \|x_k^{(1)} - \hat{x}_k^{(1)}\|_2 &= \|x_k - p_k \Gamma_k^{(1)} - (\hat{x}_k - \hat{p}_k \hat{\Gamma}_k^{(1)})\|_2 \leq \|x_k - \hat{x}_k\|_2 + \|p_k \Gamma_k^{(1)} - \hat{p}_k \hat{\Gamma}_k^{(1)}\|_2 \\ &\leq \|x_k - \hat{x}_k\|_2 + \|p_k (\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)})\|_2 + \|(p_k - \hat{p}_k) \hat{\Gamma}_k^{(1)}\|_2 \\ &\leq \|x_k - \hat{x}_k\|_2 + \|p_k (\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)})\|_2 + \|(p_k - \hat{p}_k) \Gamma_k^{(1)}\|_2 \\ &\quad + \|(p_k - \hat{p}_k) (\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)})\|_2 \\ &= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2), \end{aligned} \quad (103)$$

$$\begin{aligned} \|r_k^{(1)} - \hat{r}_k^{(1)}\|_2 &= \|r_k - q_k \Gamma_k^{(1)} - (\hat{r}_k - \hat{q}_k \hat{\Gamma}_k^{(1)})\|_2 \leq \|r_k - \hat{r}_k\|_2 + \|q_k \Gamma_k^{(1)} - \hat{q}_k \hat{\Gamma}_k^{(1)}\|_2 \\ &\leq \|r_k - \hat{r}_k\|_2 + \|q_k (\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)})\|_2 + \|(q_k - \hat{q}_k) \hat{\Gamma}_k^{(1)}\|_2 \\ &\leq \|r_k - \hat{r}_k\|_2 + \|q_k (\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)})\|_2 + \|(q_k - \hat{q}_k) \Gamma_k^{(1)}\|_2 \\ &\quad + \|(q_k - \hat{q}_k) (\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)})\|_2 \\ &= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2). \end{aligned} \quad (104)$$

Then, (75) holds as

$$\begin{aligned} \|x_k^{(2)} - \hat{x}_k^{(2)}\|_2 &= \|x_k^{(1)} + \beta_k r_k^{(1)} - (\hat{x}_k^{(1)} + \beta_k \hat{r}_k^{(1)})\|_2 \\ &\leq \|x_k^{(1)} - \hat{x}_k^{(1)}\|_2 + \beta_k \|r_k^{(1)} - \hat{r}_k^{(1)}\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2) \end{aligned}$$

is a consequence of (103) and (104).

Next, we prove (76), which follows a very similar procedure in proving (73). For the completeness, we still give the proof.

We have

$$|\Gamma_k^{(2)}| = \left| \beta_k \frac{q_k^T r_k^{(1)}}{p_k^T q_k} \right| \leq \frac{\beta_k \|q_k\|_2 \cdot \|r_k^{(1)}\|_2}{\underline{\kappa} \|x_{k-m_k} - x^*\|_2^2} = \frac{\mathcal{O}(\|x_{k-m_k} - x^*\|_2^2)}{\underline{\kappa} \|x_{k-m_k} - x^*\|_2^2} = \mathcal{O}(1), \quad (105)$$

due to (71) and $\|r_k^{(1)}\|_2 = \|r_k - q_k \Gamma_k^{(1)}\|_2 \leq \|r_k\|_2 + \|q_k\|_2 |\Gamma_k^{(1)}| \leq \eta \|r_{k-m_k}\|_2 + \mathcal{O}(\|x_{k-m_k} - x^*\|_2) = \mathcal{O}(\|x_{k-m_k} - x^*\|_2)$.

Next, if $q_k^T r_k^{(1)} \neq 0$, then

$$\begin{aligned} |\Gamma_k^{(2)} - \hat{\Gamma}_k^{(2)}| &= |\Gamma_k^{(2)}| \cdot \left| 1 - \frac{\hat{\Gamma}_k^{(2)}}{\Gamma_k^{(2)}} \right| = |\Gamma_k^{(2)}| \cdot \left| 1 - \frac{\hat{q}_k^T \hat{r}_k^{(1)}}{\hat{p}_k^T \hat{q}_k} \cdot \frac{p_k^T q_k}{q_k^T r_k^{(1)}} \right| = |\Gamma_k^{(2)}| \cdot \left| 1 - \frac{\hat{q}_k^T \hat{r}_k^{(1)}}{q_k^T r_k^{(1)}} \cdot \frac{p_k^T q_k}{\hat{p}_k^T \hat{q}_k} \right| \\ &= |\Gamma_k^{(2)}| \cdot |a_2(1 - b_2) + b_2| \leq |\Gamma_k^{(2)}| \cdot (|a_2| + |b_2| + |a_2 b_2|), \end{aligned} \quad (106)$$

where $a_2 := 1 - \frac{\hat{q}_k^T \hat{r}_k^{(1)}}{q_k^T r_k^{(1)}}$, $b_2 := 1 - \frac{p_k^T q_k}{\hat{p}_k^T \hat{q}_k}$. We have

$$|\Gamma_k^{(2)}| \cdot |a_2| = \left| \beta_k \frac{q_k^T r_k^{(1)}}{p_k^T q_k} \right| \cdot \left| \frac{q_k^T r_k^{(1)} - \hat{q}_k^T \hat{r}_k^{(1)}}{q_k^T r_k^{(1)}} \right| = \left| \beta_k \frac{q_k^T r_k^{(1)} - \hat{q}_k^T \hat{r}_k^{(1)}}{p_k^T q_k} \right|. \quad (107)$$

From (71) and (74) which have been proved,

$$|q_k^T(r_k^{(1)} - \hat{r}_k^{(1)})| \leq \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3).$$

We also have

$$\begin{aligned} |(q_k - \hat{q}_k)^T \hat{r}_k^{(1)}| &= |(q_k - \hat{q}_k)^T A(\hat{x}_k^{(1)} - x^*)| \\ &= |(q_k - \hat{q}_k)^T A(x_k^{(1)} - x^*) - (q_k - \hat{q}_k)^T A(x_k^{(1)} - \hat{x}_k^{(1)})| \\ &\leq |(q_k - \hat{q}_k)^T A(x_k - p_k \Gamma_k^{(1)} - x^*)| + |(q_k - \hat{q}_k)^T A(x_k^{(1)} - \hat{x}_k^{(1)})| \\ &= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3) + \hat{\kappa}^2 \mathcal{O}(\|x_{k-m_k} - x^*\|_2^4) = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3). \end{aligned}$$

Hence,

$$\begin{aligned} |q_k^T r_k^{(1)} - \hat{q}_k^T \hat{r}_k^{(1)}| &= |q_k^T(r_k^{(1)} - \hat{r}_k^{(1)}) + (q_k - \hat{q}_k)^T \hat{r}_k^{(1)}| \\ &\leq |q_k^T(r_k^{(1)} - \hat{r}_k^{(1)})| + |(q_k - \hat{q}_k)^T \hat{r}_k^{(1)}| + |(q_k - \hat{q}_k)^T(r_k^{(1)} - \hat{r}_k^{(1)})| \\ &= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3). \end{aligned} \quad (108)$$

Combining (107) and (108), we have

$$|\Gamma_k^{(2)}| \cdot |a_2| \leq \frac{\hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3)}{\underline{\kappa} \|x_{k-m_k} - x^*\|_2^2} = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2). \quad (109)$$

From (101),

$$|b_2| = |b_1| = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2). \quad (110)$$

As a result, by (109), (110), (105), and (106), we have

$$|\Gamma_k^{(2)} - \hat{\Gamma}_k^{(2)}| \leq |\Gamma_k^{(2)}| |a_2| + |\Gamma_k^{(2)}| |b_2| + |\Gamma_k^{(2)}| |a_2| |b_2| = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2). \quad (111)$$

Now consider the case that $q_k^T r_k^{(1)} = 0$. It is clear that $\Gamma_k^{(2)} = 0$. Then

$$|\Gamma_k^{(2)} - \hat{\Gamma}_k^{(2)}| = \left| \frac{\hat{q}_k^T \hat{r}_k^{(1)}}{\hat{p}_k^T \hat{q}_k} \right| \leq \frac{\hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^3)}{\frac{1}{2} \underline{\kappa} \|x_{k-m_k} - x^*\|_2^2} = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2),$$

where the inequality is due to (108) and (100).

Hence the (76) holds.

Finally, we prove $\|x_{k+1} - \hat{x}_{k+1}\|_2 = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2)$. It follows from (8c) that

$$\begin{aligned} \|x_{k+1} - \hat{x}_{k+1}\|_2 &= \|x_k^{(2)} - p_k \Gamma_k^{(2)} - (\hat{x}_k^{(2)} - \hat{p}_k \hat{\Gamma}_k^{(2)})\|_2 \\ &\leq \|x_k^{(2)} - \hat{x}_k^{(2)}\|_2 + \|p_k \Gamma_k^{(2)} - \hat{p}_k \hat{\Gamma}_k^{(2)}\|_2 \\ &\leq \|x_k^{(2)} - \hat{x}_k^{(2)}\|_2 + \|(p_k - \hat{p}_k) \Gamma_k^{(2)}\|_2 + \|p_k (\Gamma_k^{(2)} - \hat{\Gamma}_k^{(2)})\|_2 \\ &\quad + \|(\hat{p}_k - p_k) (\Gamma_k^{(2)} - \hat{\Gamma}_k^{(2)})\|_2 \\ &= \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2), \end{aligned}$$

where the last equality is due to (75), (71), (72), and (76) that have been proved. Therefore, the relation (68) holds.

As a result, we complete the induction. The relations (67) and (68) are proved. \square

Based on Lemma 1 and Corollary 1, we can obtain the convergence theorem of the restarted Min-AM.

Theorem 6. *Suppose that Assumption 1 holds for the optimization problem (1). Let $\{x_k\}$ denote the sequence of iterates generated by the restarted Min-AM, and define $\theta_k = \|I - \beta_k A\|_2$. Assume $\beta_j \in [\beta, \beta']$ ($j \geq 0$) for some positive constants β and β' . If $\|\nabla f(x_j)\|_2 \leq \eta_0 \|\nabla f(x_0)\|_2$ ($0 \leq j \leq k$) for a constant $\eta_0 > 0$, and x_0 is sufficiently close to x^* , then*

$$\|x_{k+1} - x^*\|_A \leq \theta_k \min_{p \in \mathcal{P}_{m_k}, p(0)=1} \|p(A)(x_{k-m_k} - x^*)\|_A + \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2), \quad (112)$$

where \mathcal{P}_{m_k} is the space of polynomials of degree not exceeding m_k . As a result,

$$\|x_{k+1} - x^*\|_A \leq 2\theta_k \left(\frac{\sqrt{L/\mu} - 1}{\sqrt{L/\mu} + 1} \right)^{m_k} \|x_{k-m_k} - x^*\|_A + \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2). \quad (113)$$

Proof. From Corollary 1, the auxiliary sequence $\{\hat{x}_k\}$ in Definition 1 satisfies

$$\|\hat{x}_{k+1} - x^*\|_A \leq \theta_k \min_{p \in \mathcal{P}_{m_k}, p(0)=1} \|p(A)(x_{k-m_k} - x^*)\|_A. \quad (114)$$

Then (112) follows from (114), (68), and the fact that $\|x_{k+1} - \hat{x}_{k+1}\|_A \leq \sqrt{L} \|x_{k+1} - \hat{x}_{k+1}\|_2$.

Since A is SPD, we can choose the Chebyshev polynomial for $p \in \mathcal{P}_{m_k}$ to obtain (e.g., see Section 6.11.3 in [53])

$$\min_{p \in \mathcal{P}_{m_k}, p(0)=1} \|p(A)\|_2 \leq \min_{p \in \mathcal{P}_{m_k}, p(0)=1} \max_{\lambda \in [\mu, L]} |p(\lambda)| \leq 2 \left(\frac{\sqrt{L/\mu} - 1}{\sqrt{L/\mu} + 1} \right)^{m_k}. \quad (115)$$

Like (63) in the proof of Corollary 1, we also have $\|p(A)\|_A = \|p(A)\|_2$. Then $\|p(A)(x_{k-m_k} - x^*)\|_A \leq \|p(A)\|_A \|x_{k-m_k} - x^*\|_A = \|p(A)\|_2 \|x_{k-m_k} - x^*\|_A$. With (115) and (112), the bound (113) holds. \square

Remark 8. In Theorem 6, we do not assume each β_k is chosen such that $\theta_k = \|I - \beta_k A\|_2 < 1$. It is expected that when θ_k is not too large, and m_k is sufficiently large, the minimization problem on the right-hand side of (112) can dominate the convergence rate, thus leading to fast convergence. However, in the case that an improper choice of β_k causes $\theta_k \gg 1$, the restarted Min-AM may have an erratic convergence behaviour or even diverge if m_k is too small.

Remark 9. From (112), a large m_k can make the first-order term diminish significantly. However, the high-order terms of errors $\hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_2^2)$ can be large since they are accumulated from the last restart. In practice, we can first choose large m , small τ , and large η in (11a)-(11c) so that the restarts do not occur too frequently. When the restarted Min-AM has a problematic convergence behaviour due to the high nonlinearity of ∇f , more frequent restarts are necessary.

Now, we give the proof of Theorem 2, which can be obtained from Theorem 6.

Proof of Theorem 2. Under the assumptions of Theorem 2, we prove that

$$\|x_j - x^*\|_A \leq \|x_0 - x^*\|_A, \quad j = 0, \dots, k, \quad (116)$$

and there exists a constant $\eta_0 > 0$ such that

$$\|\nabla f(x_j)\|_2 \leq \eta_0 \|\nabla f(x_0)\|_2, \quad j = 0, \dots, k. \quad (117)$$

We prove (116) and (117) by induction. For $k = 0$, (116) and (117) hold. Suppose that $k \geq 0$, and the results hold for k . We establish the results for $k + 1$. From (114) in the proof of Theorem 6, it follows that

$$\|\hat{x}_{k+1} - x^*\|_A \leq \theta_k \|x_{k-m_k} - x^*\|_A \leq \theta \|x_{k-m_k} - x^*\|_A. \quad (118)$$

Since $\theta_k \leq \theta$, we know $\beta_k \in [\frac{1-\theta}{\mu}, \frac{1+\theta}{L}]$. With the inductive hypothesis (117), for sufficiently small $\|x_0 - x^*\|_2$, the relations (67) and (68) hold. From (68), it follows that $\|x_{k+1} - \hat{x}_{k+1}\|_A = \hat{\kappa} \mathcal{O}(\|x_{k-m_k} - x^*\|_A^2)$. Then, there exists a constant c_1 such that $\|x_{k+1} - \hat{x}_{k+1}\|_A \leq \hat{\kappa} c_1 \|x_{k-m_k} - x^*\|_A^2$. With (118), we have

$$\|x_{k+1} - x^*\|_A \leq \theta \|x_{k-m_k} - x^*\|_A + \hat{\kappa} c_1 \|x_{k-m_k} - x^*\|_A^2.$$

Hence, by choosing $\|x_0 - x^*\|_2 \leq \frac{1-\theta}{2\sqrt{L}\hat{\kappa}c_1}$, which ensures $\|x_{k-m_k} - x^*\|_A \leq \|x_0 - x^*\|_A \leq \sqrt{L}\|x_0 - x^*\|_2 \leq \frac{1-\theta}{2\hat{\kappa}c_1}$ due to the inductive hypothesis (116), it follows that

$$\|x_{k+1} - x^*\|_A \leq \frac{1+\theta}{2}\|x_{k-m_k} - x^*\|_A < \|x_{k-m_k} - x^*\|_A \leq \|x_0 - x^*\|_A,$$

namely (116) for $k+1$. Also, $\|x_{k+1} - x^*\|_2 \leq \frac{1}{\sqrt{\mu}}\|x_{k+1} - x^*\|_A \leq \frac{1}{\sqrt{\mu}}\|x_0 - x^*\|_A \leq \frac{\sqrt{L}}{\sqrt{\mu}}\|x_0 - x^*\|_2$.

So we can impose $\|x_0 - x^*\|_2 \leq \frac{\sqrt{\mu}\hat{\rho}}{\sqrt{L}}$ to ensure $x_{k+1} \in \mathcal{B}_{\hat{\rho}}(x^*)$, which further yields that $\|r_{k+1}\|_2 \leq L\|x_{k+1} - x^*\|_2 \leq \frac{L\sqrt{L}}{\sqrt{\mu}}\|x_0 - x^*\|_2 \leq \frac{L\sqrt{L}}{\mu\sqrt{\mu}}\|r_0\|_2$, namely (117) for $k+1$, and $\eta_0 = \frac{L\sqrt{L}}{\mu\sqrt{\mu}}$. Hence, we complete the induction. The convergence result (13) follows from (113) in Theorem 6.

If $m_k = d$, then $\hat{x}_{k+1} = x^*$ from Corollary 1. Therefore, $\|x_{k+1} - x^*\|_2 = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2^2)$. \square

D Details of the eigenvalue estimation procedure

D.1 Analysis of the quadratic case

The procedure of the basic Min-AM leads to three-term recurrences for P_k and Q_k .

Proposition 3. *Under the same assumptions of Theorem 1 for solving strongly convex quadratic optimization (47), we have*

$$Ap_k = t_k^{(k-1)}p_{k-1} + t_k^{(k)}p_k + t_k^{(k+1)}p_{k+1}, \quad (119)$$

$$Aq_k = t_k^{(k-1)}q_{k-1} + t_k^{(k)}q_k + t_k^{(k+1)}q_{k+1}, \quad (120)$$

where the coefficients are given by

$$t_k^{(k-1)} = \frac{\phi_{k-1}}{\beta_{k-1}(1 - \Gamma_k^{(1)})}, \quad t_k^{(k)} = \frac{1}{1 - \Gamma_k^{(1)}} \left(\frac{1}{\beta_{k-1}} - \frac{\phi_k}{\beta_k} \right), \quad t_k^{(k+1)} = -\frac{1}{\beta_k(1 - \Gamma_k^{(1)})}, \quad (121)$$

and $\phi_k := \Gamma_k^{(1)} + \Gamma_k^{(2)} + \zeta_{k+1}$. Then $\phi_k = \Gamma_k^{(2)}$ and there exists a tridiagonal matrix $\bar{T}_k \in \mathbb{R}^{(k+1) \times k}$, such that

$$AP_k = P_{k+1}\bar{T}_k, \quad AQ_k = Q_{k+1}\bar{T}_k. \quad (122)$$

Proof. The relations can be verified by direct computation. In the $(k+1)$ -th iteration, since

$$x_{k+1} = x_k^{(2)} - p_k\Gamma_k^{(2)} = x_k^{(1)} + \beta_k r_k^{(1)} - p_k\Gamma_k^{(2)} = x_k + \beta_k r_k^{(1)} - p_k(\Gamma_k^{(1)} + \Gamma_k^{(2)}), \quad (123)$$

we have

$$\begin{aligned} p_{k+1} &= \Delta x_k - p_k \zeta_{k+1} = \beta_k r_k^{(1)} - p_k(\Gamma_k^{(1)} + \Gamma_k^{(2)} + \zeta_{k+1}) \\ &= \beta_k(r_k - q_k\Gamma_k^{(1)}) - p_k\phi_k = \beta_k(r_{k-1}^{(2)} - q_{k-1}\Gamma_{k-1}^{(2)} - q_k\Gamma_k^{(1)}) - p_k\phi_k \\ &= \beta_k(I - \beta_{k-1}A)r_{k-1}^{(1)} - \beta_k q_{k-1}\Gamma_{k-1}^{(2)} - \beta_k q_k\Gamma_k^{(1)} - p_k\phi_k \\ &= \beta_k(I - \beta_{k-1}A)\frac{p_k + p_{k-1}\phi_{k-1}}{\beta_{k-1}} - \beta_k q_{k-1}\Gamma_{k-1}^{(2)} - \beta_k q_k\Gamma_k^{(1)} - p_k\phi_k \\ &= \frac{\beta_k}{\beta_{k-1}}(p_k + p_{k-1}\phi_{k-1}) - \beta_k A(p_k + p_{k-1}\phi_{k-1}) + \beta_k Ap_{k-1}\Gamma_{k-1}^{(2)} + \beta_k Ap_k\Gamma_k^{(1)} - p_k\phi_k \\ &= \left(\frac{\beta_k}{\beta_{k-1}} - \phi_k \right) p_k + \left(\frac{\beta_k}{\beta_{k-1}}\phi_{k-1} \right) p_{k-1} + (\beta_k\Gamma_k^{(1)} - \beta_k)Ap_k \\ &\quad + (\beta_k\Gamma_{k-1}^{(2)} - \beta_k\phi_{k-1})Ap_{k-1} \\ &= \left(\frac{\beta_k}{\beta_{k-1}} - \phi_k \right) p_k + \left(\frac{\beta_k}{\beta_{k-1}}\phi_{k-1} \right) p_{k-1} + (\beta_k\Gamma_k^{(1)} - \beta_k)Ap_k, \end{aligned} \quad (124)$$

where the last term vanishes due to

$$\begin{aligned} \Gamma_{k-1}^{(2)} - \phi_{k-1} &= \Gamma_{k-1}^{(2)} - \Gamma_{k-1}^{(1)} - \Gamma_{k-1}^{(2)} - \zeta_k \\ &= -\frac{p_{k-1}^T r_{k-1}}{p_{k-1}^T q_{k-1}} - \frac{p_{k-1}^T \Delta r_{k-1}}{p_{k-1}^T q_{k-1}} = -\frac{p_{k-1}^T r_k}{p_{k-1}^T q_{k-1}} = 0, \end{aligned} \quad (125)$$

since $r_k \perp p_{k-1}$. It also indicates that

$$\phi_k = \Gamma_k^{(1)} + \Gamma_k^{(2)} + \zeta_{k+1} = \Gamma_k^{(2)} + \frac{p_k^T r_k}{p_k^T q_k} + \frac{p_k^T \Delta r_k}{p_k^T q_k} = \Gamma_k^{(2)} + \frac{p_k^T r_{k+1}}{p_k^T q_k} = \Gamma_k^{(2)}. \quad (126)$$

Now, it follows from Equation (124) that

$$Ap_k = \frac{1}{1 - \Gamma_k^{(1)}} \frac{\phi_{k-1}}{\beta_{k-1}} p_{k-1} + \frac{1}{1 - \Gamma_k^{(1)}} \left(\frac{1}{\beta_{k-1}} - \frac{\phi_k}{\beta_k} \right) p_k + \frac{-1}{\beta_k(1 - \Gamma_k^{(1)})} p_{k+1}. \quad (127)$$

The second relation about Aq_k can be obtained by noticing that $q_k = -Ap_k$. The tridiagonal matrix \bar{T}_k is

$$\bar{T}_k = \begin{pmatrix} t_1^{(1)} & t_2^{(1)} & & & & \\ t_1^{(2)} & t_2^{(2)} & t_3^{(2)} & & & \\ & t_2^{(3)} & t_3^{(3)} & t_4^{(3)} & & \\ & & \dots & & & \\ & & & t_{k-2}^{(k-1)} & t_{k-1}^{(k-1)} & t_k^{(k-1)} \\ & & & & t_{k-1}^{(k)} & t_k^{(k)} \\ & & & & & t_k^{(k+1)} \end{pmatrix} \in \mathbb{R}^{(k+1) \times k}. \quad (128)$$

□

Algorithm 3 k -step A -norm based Lanczos Algorithm for $Ax = \lambda x$, where A is SPD.

Input: Starting vector $v \in \mathbb{R}^d$

Output: Approximate eigenvalues $\{\theta_i\}_{i=1}^k$

- 1: $r = v, v_0 = 0$
 - 2: $\beta_0 = \|r\|_A$
 - 3: **for** $j = 1, 2, \dots$, until convergence **do**
 - 4: $v_j = r/\beta_{j-1}$
 - 5: $r = Av_j$
 - 6: $r = r - v_{j-1}\beta_{j-1}$
 - 7: $\alpha_j = v_j^T Ar$
 - 8: $r = r - v_j\alpha_j$
 - 9: $\beta_j = \|r\|_A$
 - 10: **if** $j = k$ **then**
 - 11: **break**
 - 12: **end if**
 - 13: **end for**
 - 14: Construct tridiagonal $T_k \in \mathbb{R}^{k \times k}$: The main diagonal is $\{\alpha_i\}_{i=1}^k$; the subdiagonals are $\{\beta_i\}_{i=1}^{k-1}$
 - 15: Compute the eigenvalue decomposition $T_k = S_k \Theta^{(k)} S_k^T$
 - 16: **return** $\Theta^{(k)}$
-

D.2 Proof of Theorem 3

For solving general nonlinear optimization, the iterations between two successive restarts in the restarted Min-AM are the basic Min-AM iterations. Hence, at the $(k+1)$ -th iteration, a tridiagonal matrix $T_k \in \mathbb{R}^{m_k \times m_k}$ can be constructed based on the coefficients starting from the $(k - m_k)$ -th iteration. We formulate it in the following definition.

Definition 2. For solving optimization problem (1) with the restarted Min-AM (see Algorithm 1), suppose that $m_k \geq 1$, and $t_k^{(k-1)}, t_k^{(k)}, t_k^{(k+1)}$ are obtained by (121). Then $T_k \in \mathbb{R}^{m_k \times m_k}$ is defined as $T_k = t_k$ if $m_k = 1$, and $T_k = (\bar{T}_{k-1}, t_k)$ if $m_k \geq 2$. Here, $\bar{T}_k = (T_k^T, t_k^{(k+1)} e_{m_k})^T \in \mathbb{R}^{(m_k+1) \times m_k}$, where e_{m_k} is the last column of I_{m_k} ; $t_k = t_k^{(k)}$ if $m_k = 1$, and $t_k = (0, \dots, 0, t_k^{(k-1)}, t_k^{(k)})^T \in \mathbb{R}^{m_k}$ if $m_k \geq 2$.

Now, we prove Theorem 3.

Proof of Theorem 3. Without loss of generality, we consider the case that no restart happens during the iterations, i.e. $m_k = k$. Let $A := \nabla^2 f(x^*)$.

The A -norm based Lanczos algorithm with a starting vector v for computing eigenvalues of an SPD matrix A is described in Algorithm 3, which is a modification of the Algorithm 4.6 in [6] with A -norm. It forms an A -orthonormal basis V_k of the Krylov subspace $\mathcal{K}_k(A, \hat{r}_0)$ by the Lanczos A -orthogonalization, i.e. $V_k^T A V_k = I$, where $\hat{r}_0 = A(x^* - x_0)$. Then the A -norm based Lanczos algorithm seeks to find $\tilde{\lambda} \in \mathbb{R}$ and $z \in \mathbb{R}^k$ such that

$$(A - \tilde{\lambda}I)V_k z \perp A V_k.$$

Then we have

$$V_k^T A A V_k z = (A V_k)^T A V_k z = \tilde{\lambda} (A V_k)^T V_k z = \tilde{\lambda} z. \quad (129)$$

The eigenvalues of $(A V_k)^T A V_k$ are the approximations to the true eigenvalues of A .

For Min-AM, we still use an auxiliary solving procedure of Min-AM on the local quadratic approximation (65), i.e. Process II in Definition 1. The same as $P_k, T_k, \bar{T}_k, t_k^{(k-1)}, t_k^{(k)}, t_k^{(k+1)}$ in Process I, the notations $\hat{P}_k, \hat{T}_k, \hat{\bar{T}}_k, \hat{t}_k^{(k-1)}, \hat{t}_k^{(k)}, \hat{t}_k^{(k+1)}$ are defined for Process II, correspondingly. Then the Min-AM in Process II seeks $\hat{\lambda} \in \mathbb{R}$ and $y \in \mathbb{R}^k$, such that

$$(A - \hat{\lambda}I)\hat{P}_k y \perp A \hat{P}_k,$$

namely,

$$(A \hat{P}_k)^T A \hat{P}_k y = \hat{\lambda} (A \hat{P}_k)^T \hat{P}_k y.$$

According to Proposition 3, we have $A \hat{P}_k = \hat{P}_{k+1} \hat{\bar{T}}_k$. Then

$$(A \hat{P}_k)^T A \hat{P}_k = \hat{P}_k^T A \hat{P}_{k+1} \hat{\bar{T}}_k = \hat{P}_k^T A \hat{P}_k \hat{T}_k.$$

Then the eigenvalues of

$$\hat{T}_k = (\hat{P}_k^T A \hat{P}_k)^{-1} (A \hat{P}_k)^T A \hat{P}_k$$

are the approximations of the true eigenvalues of A . \hat{T}_k is tridiagonal but generally not symmetric. However, it is clear that \hat{T}_k is similar to

$$(\hat{P}_k^T A \hat{P}_k)^{1/2} \hat{T}_k (\hat{P}_k^T A \hat{P}_k)^{-1/2} = (\hat{P}_k^T A \hat{P}_k)^{-1/2} \hat{P}_k^T A A \hat{P}_k (\hat{P}_k^T A \hat{P}_k)^{-1/2}, \quad (130)$$

which is symmetric. The columns of $U_k := \hat{P}_k (\hat{P}_k^T A \hat{P}_k)^{-1/2}$ are A -orthonormal and also span the same Krylov subspace $\mathcal{K}_k(A, \hat{r}_0)$, which implies that there exists an orthonormal matrix $S \in \mathbb{R}^{k \times k}$, such that $V_k = U_k S$. Therefore,

$$V_k^T A A V_k = S^T U_k^T A A U_k S. \quad (131)$$

Since $S^T S = I$, the eigenvalues of (129) and (130) are identical.

Recall that we have the explicit forms of T_k, \hat{T}_k by (121), and we restate it here for convenience.

$$t_k^{(k-1)} = \frac{\phi_{k-1}}{\beta_{k-1}(1 - \Gamma_k^{(1)})}, \quad t_k^{(k)} = \frac{1}{1 - \Gamma_k^{(1)}} \left(\frac{1}{\beta_{k-1}} - \frac{\phi_k}{\beta_k} \right), \quad t_k^{(k+1)} = -\frac{1}{\beta_k(1 - \Gamma_k^{(1)})}, \quad (132a)$$

$$\hat{t}_k^{(k-1)} = \frac{\hat{\phi}_{k-1}}{\beta_{k-1}(1 - \hat{\Gamma}_k^{(1)})}, \quad \hat{t}_k^{(k)} = \frac{1}{1 - \hat{\Gamma}_k^{(1)}} \left(\frac{1}{\beta_{k-1}} - \frac{\hat{\phi}_k}{\beta_k} \right), \quad \hat{t}_k^{(k+1)} = -\frac{1}{\beta_k(1 - \hat{\Gamma}_k^{(1)})}, \quad (132b)$$

and $\phi_k := \Gamma_k^{(1)} + \Gamma_k^{(2)} + \zeta_{k+1}$, $\hat{\phi}_k := \hat{\Gamma}_k^{(1)} + \hat{\Gamma}_k^{(2)} + \hat{\zeta}_{k+1}$. With the already proved relations (69)-(76), we can compute the bounds of $t_k^{(j)} - \hat{t}_k^{(j)}$, $j = k-1, k, k+1$. Let $\epsilon := \hat{\kappa} \mathcal{O}(\|x_0 - x^*\|_2)$ to simplify the notation. Then

$$\frac{1}{1 - \Gamma_k^{(1)}} - \frac{1}{1 - \hat{\Gamma}_k^{(1)}} = \frac{1}{1 - \Gamma_k^{(1)}} - \frac{1}{1 - \Gamma_k^{(1)} + \epsilon} = \frac{\epsilon}{(1 - \Gamma_k^{(1)})(1 - \Gamma_k^{(1)} + \epsilon)} = \hat{\kappa} \mathcal{O}(\|x_0 - x^*\|_2), \quad (133)$$

where we use the fact that $|\Gamma_k^{(1)} - \hat{\Gamma}_k^{(1)}| = \hat{\kappa}\mathcal{O}(\|x_0 - x^*\|_2)$, and the assumption that $|1 - \Gamma_k^{(1)}| \geq \tau_0$. We also have $|\phi_k| = \mathcal{O}(1)$ and $|\hat{\phi}_k - \phi_k| = \hat{\kappa}\mathcal{O}(\|x_{k-m_k} - x^*\|_2)$. Hence,

$$|t_k^{(k-1)} - \hat{t}_k^{(k-1)}| \leq \left| \frac{\phi_{k-1} - \hat{\phi}_{k-1}}{\beta_{k-1}(1 - \Gamma_k^{(1)})} \right| + \left| \frac{\hat{\phi}_{k-1}}{\beta_{k-1}} \left(\frac{1}{1 - \Gamma_k^{(1)}} - \frac{1}{1 - \hat{\Gamma}_k^{(1)}} \right) \right| = \hat{\kappa}\mathcal{O}(\|x_0 - x^*\|_2), \quad (134)$$

$$\begin{aligned} |t_k^{(k)} - \hat{t}_k^{(k)}| &\leq \left| \frac{1}{\beta_{k-1}} \left(\frac{1}{1 - \Gamma_k^{(1)}} - \frac{1}{1 - \hat{\Gamma}_k^{(1)}} \right) \right| + \left| \frac{\phi_k - \hat{\phi}_k}{\beta_k(1 - \Gamma_k^{(1)})} \right| \\ &\quad + \left| \frac{\hat{\phi}_k}{\beta_k} \left(\frac{1}{1 - \Gamma_k^{(1)}} - \frac{1}{1 - \hat{\Gamma}_k^{(1)}} \right) \right| \\ &= \hat{\kappa}\mathcal{O}(\|x_0 - x^*\|_2), \end{aligned} \quad (135)$$

$$|t_k^{(k+1)} - \hat{t}_k^{(k+1)}| = \hat{\kappa}\mathcal{O}(\|x_0 - x^*\|_2). \quad (136)$$

With $m_k \leq m$, we have that

$$\|T_k - \hat{T}_k\|_2 = \hat{\kappa}\mathcal{O}(\|x_0 - x^*\|_2). \quad (137)$$

Define $\hat{D}_k = \hat{P}_k^T A \hat{P}_k$. It is clear that \hat{D}_k is a diagonal matrix, and the i -th diagonal element is $-\hat{p}_i^T \hat{q}_i$. From (100), (71), and (72), it follows that

$$\|\hat{D}_k\|_2 \|\hat{D}_k^{-1}\|_2 = \frac{\max_i \{|\hat{p}_i^T \hat{q}_i|\}}{\min_j \{|\hat{p}_j^T \hat{q}_j|\}} = \mathcal{O}(1). \quad (138)$$

By (130), we know that $\hat{D}_k^{1/2} \hat{T}_k \hat{D}_k^{-1/2}$ can be diagonalized. There exists orthonormal matrix $W_k \in \mathbb{R}^{k \times k}$ and diagonal matrix $\wedge_k \in \mathbb{R}^{k \times k}$, such that

$$\hat{D}_k^{1/2} \hat{T}_k \hat{D}_k^{-1/2} = W_k \wedge_k W_k^T. \quad (139)$$

Therefore, $\hat{T}_k = \hat{D}_k^{-1/2} W_k \wedge_k W_k^T \hat{D}_k^{1/2}$ can be diagonalized. Then with (137), (138), and applying the Bauer-Fike Theorem [33], we know that for an eigenvalue λ of T_k , the following result holds:

$$\begin{aligned} \min_{\hat{\lambda} \in \sigma(\hat{T}_k)} |\hat{\lambda} - \lambda| &\leq \|\hat{D}_k^{-1/2} W_k\|_2 \|W_k^T \hat{D}_k^{1/2}\|_2 \|T_k - \hat{T}_k\|_2 \\ &\leq \|\hat{D}_k^{-1/2}\|_2 \|\hat{D}_k^{1/2}\|_2 \|T_k - \hat{T}_k\|_2 = \hat{\kappa}\mathcal{O}(\|x_0 - x^*\|_2), \end{aligned} \quad (140)$$

where $\sigma(\hat{T}_k)$ denotes the spectrum of \hat{T}_k . \square

E Details of the stochastic Min-AM

Now, we prove the theorems about the stochastic Min-AM (sMin-AM) in Section 3.4. The algorithm is given in Algorithm 4. For brevity, we use δ_k to denote $\delta_k^{(2)}$, i.e. $\delta_k \equiv \delta_k^{(2)}$. Our proofs follow those of SAM [62] and ST-AM [63].

From Assumption 3, for the mini-batch gradient $f_{S_k}(x_k) = \frac{1}{n_k} \sum_{i \in S_k} f_i(x_k)$, where $n_k = |S_k|$, we have

$$\mathbb{E}[\nabla f_{S_k}(x)|x_k] = \nabla f(x_k), \quad (141a)$$

$$\mathbb{E}[\|\nabla f_{S_k}(x_k) - \nabla f(x_k)\|_2^2 | x_k] \leq \frac{\sigma^2}{n_k}. \quad (141b)$$

Note that the update of sMin-AM (Lines 10-12 in Algorithm 4) can be written as $x_{k+1} = x_k + H_k r_k$, where $r_k = -\nabla f_{S_k}(x_k)$, and for $k \geq 0$,

$$H_k = (1 - \alpha_k)\beta_k I + \alpha_k H_k^A = \beta_k I + \alpha_k (H_k^A - \beta_k I). \quad (142)$$

To prove the theorems, we need H_k to be positive definite, and show that the noise in gradient estimates is suppressed during iterations.

We first state some useful results about sMin-AM, then we prove Theorem 4 and Theorem 5.

Algorithm 4 Stochastic Min-AM

Input: $x_0 \in \mathbb{R}^d, \beta_k > 0, \alpha_k \in [0, 1], \delta_k^{(1)} \geq 0, \delta_k^{(2)} \geq 0$

- 1: $p_0, q_0 = \mathbf{0} \in \mathbb{R}^d$
- 2: **for** $k = 0, 1, \dots$, until convergence, **do**
- 3: $r_k = -\nabla f_{S_k}(x_k)$
- 4: **if** $k > 0$ **then**
- 5: $p = x_k - x_{k-1}, q = r_k - r_{k-1}$
- 6: $\zeta_k = \Phi(p_{k-1}, q_{k-1}, \delta_k^{(1)})^\dagger p_{k-1}^\top q$
- 7: $q_k = q - q_{k-1} \zeta_k, p_k = p - p_{k-1} \zeta_k$
- 8: **end if**
- 9: $\rho_k = \Phi(p_k, q_k, \delta_k^{(2)})^\dagger$
- 10: $x_k^A = x_k + (-\rho_k p_k p_k^\top + \beta_k (I - \rho_k p_k q_k^\top)(I - \rho_k q_k p_k^\top)) r_k$
- 11: $x_k^G = x_k + \beta_k r_k$
- 12: $x_{k+1} = (1 - \alpha_k) x_k^G + \alpha_k x_k^A$
- 13: Apply learning rate schedule of α_k, β_k
- 14: **end for**
- 15: **return** x_k

E.1 Some useful results

Lemma 2. Suppose that the sequence $\{x_k\}$ is generated by sMin-AM. If there are constants $\mu \in (0, 1), C_2 > 0$ such that $\alpha_k \in [0, 1 - \mu], \rho_k \leq 0$ and $-\rho_k \|p_k\|_2^2 - 2\beta_k \rho_k \|p_k\|_2 \|q_k\|_2 + \beta_k \rho_k^2 \|p_k\|_2^2 \|q_k\|_2^2 \leq \beta_k C_2$, then we have that

$$\|H_k^A - \beta_k I\|_2 \leq C_2 \beta_k, \quad (143)$$

$$\|H_k\|_2 \leq \beta_k (1 + C_2), \quad (144)$$

$$H_k \succeq \beta_k \mu I. \quad (145)$$

Proof. Since

$$\begin{aligned} H_k^A &= -\rho_k p_k p_k^\top + \beta_k (1 - \rho_k p_k q_k^\top)(1 - \rho_k q_k p_k^\top) \\ &= -\rho_k p_k p_k^\top + \beta_k I - \beta_k \rho_k q_k p_k^\top - \beta_k \rho_k p_k q_k^\top + \beta_k \rho_k^2 p_k q_k^\top q_k p_k^\top, \end{aligned}$$

we have

$$H_k^A - \beta_k I = -\rho_k p_k p_k^\top - \beta_k \rho_k q_k p_k^\top - \beta_k \rho_k p_k q_k^\top + \beta_k \rho_k^2 p_k q_k^\top q_k p_k^\top. \quad (146)$$

Hence,

$$\begin{aligned} \|H_k^A - \beta_k I\|_2 &\leq |\rho_k| \|p_k p_k^\top\|_2 + \beta_k |\rho_k| \|q_k p_k^\top\|_2 + \beta_k |\rho_k| \|p_k q_k^\top\|_2 + \beta_k \rho_k^2 \|p_k q_k^\top q_k p_k^\top\|_2 \\ &\leq |\rho_k| \|p_k\|_2^2 + \beta_k |\rho_k| \|q_k\|_2 \|p_k\|_2 + \beta_k |\rho_k| \|p_k\|_2 \|q_k\|_2 + \beta_k \rho_k^2 \|q_k\|_2^2 \|p_k\|_2^2 \\ &\leq \beta_k C_2. \end{aligned} \quad (147)$$

It follows that

$$\|H_k\|_2 \leq \|\beta_k I\|_2 + \alpha_k \|H_k^A - \beta_k I\|_2 \leq \beta_k + C_2 \beta_k = (1 + C_2) \beta_k. \quad (148)$$

Because $\rho_k \leq 0$, it follows that $H_k^A \succeq 0$. Hence $H_k \succeq (1 - \alpha_k) \beta_k I \succeq \mu \beta_k I$. \square

Remark 10. The conditions that $\rho_k \leq 0$ and $-\rho_k \|p_k\|_2^2 - 2\beta_k \rho_k \|p_k\|_2 \|q_k\|_2 + \beta_k \rho_k^2 \|p_k\|_2^2 \|q_k\|_2^2 \leq \beta_k C_2$ for a positive constant C_2 can be fulfilled with some proper choice of δ_k . For example, if C_2 satisfies $(1 + \frac{1}{2}\beta_k)^2 \leq C_2$ and we choose $\delta_k \geq \frac{1}{2} + \beta_k^{-1}$, then it can be proved that the conditions hold. Moreover, in practice, δ_k can be chosen as a constant, as shown in our experiments.

With Lemma 2, we can prove the convergence of sMin-AM for full-batch training.

Theorem 7. Suppose that Assumption 2 holds and $\{x_k\}$ is the sequence generated by full-batch sMin-AM, i.e. $n_k = T$. Given constants $\mu \in (0, 1), C_2 > 0$, let $\beta_k = \beta \in (0, \frac{\mu}{L(1+C_2)^2}]$ be a

constant, $\alpha_k \in [0, 1 - \mu]$. The $\delta_k^{(2)}$ is chosen such that $\rho_k \leq 0$ and $-\rho_k \|p_k\|_2^2 - 2\beta_k \rho_k \|p_k\|_2 \|q_k\|_2 + \beta_k \rho_k^2 \|p_k\|_2^2 \|q_k\|_2^2 \leq \beta_k C_2$. Then

$$\frac{1}{N} \sum_{k=0}^{N-1} \|\nabla f(x_k)\|_2^2 \leq \frac{2(f(x_0) - f^{low})}{N\mu\beta}, \quad (149)$$

in the N -th iterations. To ensure $\frac{1}{N} \sum_{k=0}^{N-1} \|\nabla f(x_k)\|_2^2 < \epsilon$, the number of iterations is $\mathcal{O}(1/\epsilon)$.

Proof of Theorem 7. With Lemma 2, we have

$$\|H_k r_k\|_2^2 \leq \beta^2 (1 + C_2)^2 \|r_k\|_2^2,$$

and

$$r_k^T H_k r_k \geq \beta\mu \|r_k\|_2^2.$$

Then, under Assumption 2, we have

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|_2^2 \\ &= f(x_k) - r_k^T H_k r_k + \frac{L}{2} \|H_k r_k\|_2^2 \\ &\leq f(x_k) - \beta\mu \|r_k\|_2^2 + \frac{L}{2} \beta^2 (1 + C_2)^2 \|r_k\|_2^2 \\ &= f(x_k) - \beta \left(\mu - \frac{L}{2} \beta (1 + C_2)^2 \right) \|r_k\|_2^2 \\ &\leq f(x_k) - \frac{1}{2} \beta\mu \cdot \|\nabla f(x_k)\|_2^2, \end{aligned} \quad (150)$$

where the last inequality is due to $0 < \beta \leq \frac{\mu}{L(1+C_2)^2}$. Thus, $f(x_{k+1}) - f(x_k) \leq -\frac{1}{2} \beta\mu \|\nabla f(x_k)\|_2^2$.

Summing both sides of this inequality for $k \in \{0, \dots, N-1\}$ and recalling $f(x) > f^{low}$ in Assumption 2 gives

$$f^{low} - f(x_0) \leq f(x_N) - f(x_0) \leq -\frac{1}{2} \beta\mu \sum_{k=0}^{N-1} \|\nabla f(x_k)\|_2^2.$$

Rearranging and dividing further by N yields (149). \square

The next lemmas are about the stochastic case.

Lemma 3. Suppose that Assumption 3 and Assumption 4 hold and $\{x_k\}$ is the sequence generated by sMin-AM with $\beta_k > 0$. Then

$$\mathbb{E}_{S_k} [\|H_k r_k\|_2^2] \leq \beta_k^2 (1 + C_2)^2 \cdot \left(\|\nabla f(x_k)\|_2^2 + \frac{\sigma^2}{n_k} \right), \quad (151a)$$

$$\nabla f(x_k)^T \mathbb{E}_{S_k} [H_k r_k] \leq -\frac{1}{2} \beta_k \mu \|\nabla f(x_k)\|_2^2 + \frac{1}{2} \beta_k^2 \cdot \mu^{-1} C_1^2 C_2^2 \frac{\sigma^2}{n_k}. \quad (151b)$$

Proof. (i) From Lemma 2, we have

$$\mathbb{E}_{S_k} [\|H_k r_k\|_2^2] \leq \beta_k^2 (1 + C_2)^2 \mathbb{E}_{S_k} [\|r_k\|_2^2]. \quad (152)$$

From Assumption 3, we have

$$\mathbb{E}_{S_k} [\|r_k\|_2^2] = \mathbb{E}_{S_k} [\|r_k - \mathbb{E}_{S_k} [r_k]\|_2^2] + \|\mathbb{E}_{S_k} [r_k]\|_2^2 \leq \|\nabla f(x_k)\|_2^2 + \sigma^2/n_k. \quad (153)$$

With (152) and (153), we obtain (151a).

(ii) Recalling that $H_0 = \beta_0 I$, the result holds for $k = 0$. Define $\epsilon_k = \nabla f_{S_k}(x_k) - \nabla f(x_k) = -r_k - \nabla f(x_k)$, then $H_k r_k = H_k (-\epsilon_k - \nabla f(x_k))$. First, we have

$$\nabla f(x_k)^T H_k \nabla f(x_k) \geq \beta_k \mu \|\nabla f(x_k)\|_2^2,$$

which implies

$$\mathbb{E}_{S_k}[\nabla f(x_k)^T H_k \nabla f(x_k)] \geq \beta_k \mu \|\nabla f(x_k)\|_2^2. \quad (154)$$

Let $M_k := \alpha_k(H_k^A - \beta_k I)$, then $H_k = \beta_k I + M_k$. With the assumption (141a), i.e. $\mathbb{E}_{S_k}[\epsilon_k] = 0$, we have

$$\begin{aligned} \mathbb{E}_{S_k}[\nabla f(x_k)^T H_k \epsilon_k] &= \mathbb{E}_{S_k}[\nabla f(x_k)^T (\beta_k \epsilon_k + M_k \epsilon_k)] \\ &= \beta_k \nabla f(x_k)^T \mathbb{E}_{S_k}[\epsilon_k] + \mathbb{E}_{S_k}[\nabla f(x_k)^T M_k \epsilon_k] = \mathbb{E}_{S_k}[\nabla f(x_k)^T M_k \epsilon_k]. \end{aligned}$$

Using the *Cauchy-Schwarz inequality with expectations*, we obtain

$$\begin{aligned} |\mathbb{E}_{S_k}[\nabla f(x_k)^T H_k \epsilon_k]| &= |\mathbb{E}_{S_k}[\nabla f(x_k)^T M_k \epsilon_k]| \leq \sqrt{\mathbb{E}_{S_k}[\|\nabla f(x_k)\|_2^2]} \sqrt{\mathbb{E}_{S_k}[\|M_k \epsilon_k\|_2^2]} \\ &= \|\nabla f(x_k)\|_2 \sqrt{\mathbb{E}_{S_k}[\|M_k \epsilon_k\|_2^2]}. \end{aligned} \quad (155)$$

We now bound $\|M_k \epsilon_k\|_2^2$.

$$\|M_k\|_2 = \alpha_k \|H_k^A - \beta_k I\|_2 \leq C_2 \alpha_k \beta_k. \quad (156)$$

With (156), we have $\|M_k \epsilon_k\|_2 \leq C_2 \alpha_k \beta_k \|\epsilon_k\|_2$, which implies

$$\mathbb{E}_{S_k}[\|M_k \epsilon_k\|_2^2] \leq C_2^2 \alpha_k^2 \beta_k^2 \mathbb{E}_{S_k}[\|\epsilon_k\|_2^2] \leq C_2^2 \alpha_k^2 \beta_k^2 \frac{\sigma^2}{n_k}, \quad (157)$$

where the last inequality is due to (141b). Now we can obtain the bound of $|\mathbb{E}_{S_k}[\nabla f(x_k)^T H_k \epsilon_k]|$ as follows (cf. (155)):

$$\begin{aligned} |\mathbb{E}_{S_k}[\nabla f(x_k)^T H_k \epsilon_k]| &\leq \|\nabla f(x_k)\|_2 \sqrt{\mathbb{E}_{S_k}[\|M_k \epsilon_k\|_2^2]} \\ &\leq C_2 \alpha_k \beta_k \|\nabla f(x_k)\|_2 \sqrt{\mathbb{E}_{S_k}[\|\epsilon_k\|_2^2]} \\ &\leq C_2 \alpha_k \beta_k \frac{\sigma}{\sqrt{n_k}} \|\nabla f(x_k)\|_2 \\ &= \sqrt{\beta_k \mu} \|\nabla f(x_k)\|_2 \cdot \frac{C_2 \alpha_k \beta_k}{\sqrt{\beta_k \mu}} \frac{\sigma}{\sqrt{n_k}} \\ &\leq \frac{1}{2} \beta_k \mu \|\nabla f(x_k)\|_2^2 + \frac{1}{2} \frac{C_2^2 \alpha_k^2 \beta_k^2}{\beta_k \mu} \cdot \frac{\sigma^2}{n_k}. \end{aligned} \quad (158)$$

With the inequality (154) and (158), we obtain

$$\begin{aligned} &\nabla f(x_k)^T \mathbb{E}_{S_k}[H_k r_k] \\ &= -\nabla f(x_k)^T \mathbb{E}_{S_k}[H_k (\epsilon_k + \nabla f(x_k))] \\ &= -\mathbb{E}_{S_k}[\nabla f(x_k)^T H_k \nabla f(x_k)] - \mathbb{E}_{S_k}[\nabla f(x_k)^T H_k \epsilon_k] \\ &\leq -\mathbb{E}_{S_k}[\nabla f(x_k)^T H_k \nabla f(x_k)] + |\mathbb{E}_{S_k}[\nabla f(x_k)^T H_k \epsilon_k]| \\ &\leq -\beta_k \mu \|\nabla f(x_k)\|_2^2 + \frac{1}{2} \beta_k \mu \|\nabla f(x_k)\|_2^2 + \frac{1}{2} \frac{C_2^2 \alpha_k^2 \beta_k^2}{\beta_k \mu} \cdot \frac{\sigma^2}{n_k} \\ &= -\frac{1}{2} \beta_k \mu \|\nabla f(x_k)\|_2^2 + \frac{1}{2} \frac{C_2^2 \alpha_k^2 \beta_k^2}{\beta_k \mu} \cdot \frac{\sigma^2}{n_k} \leq -\frac{1}{2} \beta_k \mu \|\nabla f(x_k)\|_2^2 + \frac{1}{2} \beta_k^2 \mu^{-1} C_1^2 C_2^2 \cdot \frac{\sigma^2}{n_k}, \end{aligned} \quad (159)$$

where the last inequality is due to $\alpha_k \leq C_1 \beta_k^{1/2}$. \square

Using Lemma 3 we obtain the descent property of sMin-AM:

Lemma 4. Suppose that Assumptions 2-4 hold, $\beta_k \in (0, \frac{\mu}{2L(1+C_2)^2}]$, and $\{x_k\}$ is the sequence generated by sMin-AM. Then

$$\mathbb{E}_{S_k}[f(x_{k+1})] \leq f(x_k) - \frac{1}{4} \beta_k \mu \|\nabla f(x_k)\|_2^2 + \frac{\beta_k^2}{2} (\mu^{-1} C_1^2 C_2^2 + L(1+C_2)^2) \frac{\sigma^2}{n_k}. \quad (160)$$

Proof. Due to Assumption 2, we have

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \nabla f(x_k)^\top (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|_2^2 \\ &= f(x_k) + \nabla f(x_k)^\top H_k r_k + \frac{L}{2} \|H_k r_k\|_2^2. \end{aligned} \quad (161)$$

Taking expectation with respect to the mini-batch S_k on both sides of (161) and using Lemma 3, we obtain

$$\begin{aligned} &\mathbb{E}_{S_k}[f(x_{k+1})] \\ &\leq f(x_k) + \nabla f(x_k)^\top \mathbb{E}_{S_k}[H_k r_k] + \frac{L}{2} \mathbb{E}_{S_k} \|H_k r_k\|_2^2 \\ &\leq f(x_k) - \frac{1}{2} \beta_k \mu \|\nabla f(x_k)\|_2^2 + \frac{1}{2} \beta_k^2 \cdot \mu^{-1} C_1^2 C_2^2 \frac{\sigma^2}{n_k} + \frac{L}{2} \beta_k^2 (1 + C_2)^2 \left(\|\nabla f(x_k)\|_2^2 + \frac{\sigma^2}{n_k} \right) \\ &= f(x_k) - \beta_k \left(\frac{1}{2} \mu - \frac{L}{2} \beta_k (1 + C_2)^2 \right) \|\nabla f(x_k)\|_2^2 + \frac{\beta_k^2}{2} (\mu^{-1} C_1^2 C_2^2 + L(1 + C_2)^2) \frac{\sigma^2}{n_k}. \end{aligned} \quad (162)$$

Then (162) combined with the assumption $\beta_k \leq \frac{\mu}{2L(1+C_2)^2}$ implies (160). \square

E.2 Proof of Theorem 4

Following the proofs in [61, 62, 63], we introduce the definition of a *supermartingale*.

Definition 3. Let $\{\mathcal{F}_k\}$ be an increasing sequence of σ -algebras. If $\{X_k\}$ is a stochastic process satisfying (i) $\mathbb{E}[|X_k|] < \infty$, (ii) $X_k \in \mathcal{F}_k$ for all k , and (iii) $\mathbb{E}[X_{k+1} | \mathcal{F}_k] \leq X_k$ for all k , then $\{X_k\}$ is called a *supermartingale*.

Proposition 4 (Supermartingale convergence theorem, see, e.g., Theorem 4.2.12 in [20]). If $\{X_k\}$ is a nonnegative supermartingale, then $\lim_{k \rightarrow \infty} X_k \rightarrow X$ almost surely and $\mathbb{E}[X] \leq \mathbb{E}[X_0]$.

Now, we prove Theorem 4 of sMin-AM.

Proof of Theorem 4. (i) Define $\phi_k = \frac{\beta_k \mu}{4} \|\nabla f(x_k)\|_2^2$ and $\tilde{L} = \frac{1}{2} (\mu^{-1} C_1^2 C_2^2 + L(1 + C_2)^2)$, $\gamma_k = f(x_k) + \tilde{L} \frac{\sigma^2}{n} \sum_{i=k}^{\infty} \beta_i^2$. Let \mathcal{F}_k be the σ -algebra measuring ϕ_k, γ_k , and x_k . From (160) we know that for any k ,

$$\begin{aligned} \mathbb{E}[\gamma_{k+1} | \mathcal{F}_k] &= \mathbb{E}[f(x_{k+1}) | \mathcal{F}_k] + \tilde{L} \frac{\sigma^2}{n} \sum_{i=k+1}^{\infty} \beta_i^2 \\ &\leq f(x_k) - \frac{1}{4} \beta_k \mu \|\nabla f(x_k)\|_2^2 + \tilde{L} \frac{\sigma^2}{n} \sum_{i=k}^{\infty} \beta_i^2 = \gamma_k - \phi_k, \end{aligned} \quad (163)$$

which implies that $\mathbb{E}[\gamma_{k+1} - f^{low} | \mathcal{F}_k] \leq \gamma_k - f^{low} - \phi_k$. Since $\phi_k \geq 0$, we have $0 \leq \mathbb{E}[\gamma_k - f^{low}] \leq \gamma_0 - f^{low} < +\infty$. As the diminishing condition (25) holds, we obtain $\mathbb{E}[f(x_k)] \leq M_f$ for some constant $M_f > 0$. According to Definition 3, $\{\gamma_k - f^{low}\}$ is a supermartingale. Therefore, Proposition 4 indicates that there exists a constant γ such that $\lim_{k \rightarrow \infty} \gamma_k = \gamma$ with probability 1, and $\mathbb{E}[\gamma] \leq \mathbb{E}[\gamma_0]$. Note that from (163) we have $\mathbb{E}[\phi_k] \leq \mathbb{E}[\gamma_k] - \mathbb{E}[\gamma_{k+1}]$. Thus,

$$\mathbb{E} \left[\sum_{k=0}^{\infty} \phi_k \right] \leq \sum_{k=0}^{\infty} (\mathbb{E}[\gamma_k] - \mathbb{E}[\gamma_{k+1}]) < +\infty,$$

which further yields that

$$\sum_{k=0}^{\infty} \phi_k = \frac{\mu}{4} \sum_{k=0}^{\infty} \beta_k \|\nabla f(x_k)\|_2^2 < +\infty \text{ with probability 1.} \quad (164)$$

Since $\sum_{k=0}^{\infty} \beta_k = +\infty$, it follows that (26) holds.

(ii) If the noisy gradient is bounded, i.e.,

$$\mathbb{E}_{\xi_k} [\|\nabla f_{\xi_k}(x_k)\|_2^2] \leq M_g, \quad (165)$$

where $M_g > 0$ is a constant, then a stronger result can be obtained.

For any given $\epsilon > 0$, according to (26), there exist infinitely many iterates x_k such that $\|\nabla f(x_k)\|_2 \leq \epsilon$. Then if (27) does not hold, there must exist two infinite sequences of indices $\{s_i\}$, $\{t_i\}$ with $t_i > s_i$, such that for $i = 0, 1, \dots$, $k = s_i + 1, \dots, t_i - 1$,

$$\|\nabla f(x_{s_i})\|_2 \geq 2\epsilon, \|\nabla f(x_{t_i})\|_2 < \epsilon, \|\nabla f(x_k)\|_2 \geq \epsilon. \quad (166)$$

Then from (164) it follows that

$$+\infty > \sum_{k=0}^{\infty} \beta_k \|\nabla f(x_k)\|_2^2 \geq \sum_{i=0}^{+\infty} \sum_{k=s_i}^{t_i-1} \beta_k \|\nabla f(x_k)\|_2^2 \geq \epsilon^2 \sum_{i=0}^{+\infty} \sum_{k=s_i}^{t_i-1} \beta_k \text{ with probability 1,}$$

which implies that

$$\sum_{k=s_i}^{t_i-1} \beta_k \rightarrow 0 \text{ with probability 1, as } i \rightarrow +\infty. \quad (167)$$

According to (144), we have

$$\begin{aligned} & \mathbb{E}[\|x_{k+1} - x_k\|_2 | x_k] \\ &= \mathbb{E}[\|H_k r_k\|_2 | x_k] \leq \beta_k (1 + C_2) \mathbb{E}[\|r_k\|_2 | x_k] \leq \beta_k (1 + C_2) (\mathbb{E}[\|r_k\|_2^2 | x_k])^{\frac{1}{2}} \leq \beta_k (1 + C_2) M_g^{\frac{1}{2}}, \end{aligned} \quad (168)$$

where the last inequalities are due to *Cauchy-Schwarz inequality* and (165). Then it follows from (168) that

$$\mathbb{E}[\|x_{t_i} - x_{s_i}\|_2] \leq (1 + C_2) M_g^{\frac{1}{2}} \sum_{k=s_i}^{t_i-1} \beta_k,$$

which together with (167) implies that $\|x_{t_i} - x_{s_i}\|_2 \rightarrow 0$ with probability 1, as $i \rightarrow +\infty$. Hence, from the Lipschitz continuity of ∇f , it follows that $\|\nabla f(x_{t_i}) - \nabla f(x_{s_i})\|_2 \rightarrow 0$ with probability 1 as $i \rightarrow +\infty$. However, this contradicts (166). Therefore, the assumption that (27) does not hold is not true. \square

E.3 Proof of Theorem 5

Proof of Theorem 5. According to (162) in Lemma 4, we have

$$\begin{aligned} & \sum_{k=0}^{N-1} \beta_k \left(\frac{1}{2} \mu - \frac{L}{2} \beta_k (1 + C_2)^2 \right) \mathbb{E} \|\nabla f(x_k)\|_2^2 \\ & \leq f(x_0) - f^{low} + \sum_{k=0}^{N-1} \frac{\beta_k^2}{2} (\mu^{-1} C_1^2 C_2^2 + L(1 + C_2)^2) \frac{\sigma^2}{n_k}, \end{aligned} \quad (169)$$

where the expectation is taken with respect to $\{S_j\}_{j=0}^{N-1}$. Define

$$P_R(k) = \text{Prob}\{R = k\} = \frac{\beta_k \left(\frac{1}{2} \mu - \frac{L}{2} \beta_k (1 + C_2)^2 \right)}{\sum_{j=0}^{N-1} \beta_j \left(\frac{1}{2} \mu - \frac{L}{2} \beta_j (1 + C_2)^2 \right)}, \quad k = 0, \dots, N-1, \quad (170)$$

then

$$\begin{aligned} \mathbb{E} [\|\nabla f(x_R)\|_2^2] &= \frac{\sum_{k=0}^{N-1} \beta_k \left(\frac{1}{2} \mu - \frac{L}{2} \beta_k (1 + C_2)^2 \right) \mathbb{E} [\|\nabla f(x_k)\|_2^2]}{\sum_{j=0}^{N-1} \beta_j \left(\frac{1}{2} \mu - \frac{L}{2} \beta_j (1 + C_2)^2 \right)} \\ &\leq \frac{D_f + \frac{\sigma^2}{2} (\mu^{-1} C_1^2 C_2^2 + L(1 + C_2)^2) \sum_{k=0}^{N-1} \beta_k^2 / n_k}{\sum_{j=0}^{N-1} \beta_j \left(\frac{1}{2} \mu - \frac{L}{2} \beta_j (1 + C_2)^2 \right)}. \end{aligned} \quad (171)$$

Let \tilde{D} be a problem-independent constant. If we choose $\beta_k = \beta := \min\{\frac{\mu}{2L(1+C_2)^2}, \frac{\tilde{D}}{\sigma\sqrt{N}}\}$, and $n_k \equiv n$, then the definition of P_R simplifies to $P_R(k) = 1/N$. From (171) we have

$$\begin{aligned}
\mathbb{E}[\|\nabla f(x_R)\|_2^2] &\leq \frac{D_f + \frac{\sigma^2}{2}(\mu^{-1}C_1^2C_2^2 + L(1+C_2)^2)\frac{N\beta^2}{n}}{\sum_{j=0}^{N-1} \beta(\frac{1}{2}\mu - \frac{\mu}{4})} \\
&= \frac{D_f + \frac{\sigma^2}{2}(\mu^{-1}C_1^2C_2^2 + L(1+C_2)^2)\frac{N\beta^2}{n}}{N\beta \cdot \frac{1}{4}\mu} \\
&= \frac{4D_f}{N\beta\mu} + \frac{\frac{\sigma^2}{2}(\mu^{-1}C_1^2C_2^2 + L(1+C_2)^2) \cdot \beta}{\frac{1}{4}n\mu} \\
&\leq \frac{4D_f}{N\mu} \max\left\{\frac{2L(1+C_2)^2}{\mu}, \frac{\sigma\sqrt{N}}{\tilde{D}}\right\} + \frac{2\sigma^2(\mu^{-1}C_1^2C_2^2 + L(1+C_2)^2)}{n\mu} \cdot \frac{\tilde{D}}{\sigma\sqrt{N}} \\
&\leq \frac{4D_f}{N\mu} \left(\frac{2L(1+C_2)^2}{\mu} + \frac{\sigma\sqrt{N}}{\tilde{D}}\right) + \frac{2\sigma(\mu^{-1}C_1^2C_2^2 + L(1+C_2)^2)\tilde{D}}{n\mu\sqrt{N}} \\
&= \frac{8D_fL(1+C_2)^2}{N\mu^2} + \frac{\sigma}{\mu\sqrt{N}} \left(\frac{4D_f}{\tilde{D}} + \frac{2(\mu^{-1}C_1^2C_2^2 + L(1+C_2)^2)\tilde{D}}{n}\right).
\end{aligned}$$

Therefore, to ensure $\mathbb{E}[\|\nabla f(x_R)\|_2^2] \leq \epsilon$, the number of iterations is $\mathcal{O}(1/\epsilon^2)$. \square

F Experimental details

We implemented the algorithms based on PyTorch² and used one GeForce RTX 2080 Ti GPU for training neural networks (except that four GPUs were used for training ResNet50 on ImageNet). We tested the basic Min-AM on solving strongly convex quadratic optimization, the restarted Min-AM on solving regularized logistic regression, and the stochastic Min-AM on training neural networks.

F.1 Strongly convex quadratic problem

The experiments on solving strongly convex quadratic optimization were conducted to verify Theorem 1 about the basic Min-AM. The problem is

$$\min_{x \in \mathbb{R}^d} f(x) := \frac{1}{2} \|Ax - b\|_2^2, \quad (172)$$

where $A \in \mathbb{R}^{\ell \times d}$, $b \in \mathbb{R}^\ell$. For the test, we first generated a random matrix $A \in \mathbb{R}^{500 \times 100}$ and a random vector $v \in \mathbb{R}^{100}$ following Gaussian distribution, then $b \in \mathbb{R}^{500}$ was obtained as $b = Av$.

The fixed step size β for the gradient descent (GD) was chosen by a grid search in $\{0.001, 0.002, \dots, 0.01\}$. We set $\beta = 0.001$ that guarantees the convergence of GD. The AM-I(1) and the full-memory AM-I (i.e. AM-I(∞)) also used the same mixing parameter setting $\beta_k = 0.001$. For Min-AM, we set the initial mixing parameter $\beta_0 = 1$, and the later mixing parameters $\{\beta_k\}$ were adaptively determined based on the eigenvalue estimates (see Section 3.3).

Figure 4(a) compares the convergence behaviours of different methods in terms of relative residual norm. It can be observed that due to the improper initial setting of β_0 for Min-AM, Min-AM does not perform well in the beginning. Nonetheless, as shown in Figure 4(b), the β_k of Min-AM can be quickly adapted to the optimal value $2/(\mu + L) = 1.67 \times 10^{-3}$ based on the eigenvalue estimates.

In Figure 4(c), we show the effects of β_k on the full-memory AM-I and Min-AM, where both methods used fixed β_k chosen from $\{0.01, 0.1, 1\}$. (For Min-AM, we disable the adaptive choice of β_k .) It should be noted that the GD method diverges when choosing the step size in $\{0.01, 0.1, 1\}$, which suggests that $\theta_k > 1$ for the Min-AM in these settings. Nevertheless, we find that Min-AM still converges to the tolerance of $\|r_k\|_2/\|r_0\|_2 \leq 10^{-11}$, which validates the convergence property (59) in Corollary 1, i.e., the minimization problem in (59) dominates the convergence when k is

² <https://pytorch.org>.

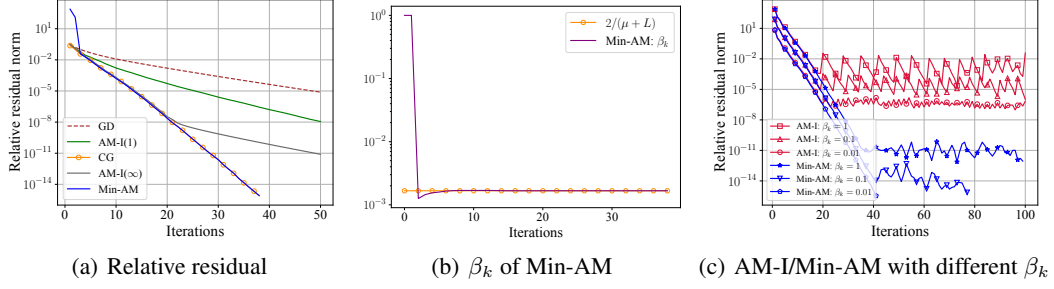


Figure 4: (a) $\|r_k\|_2/\|r_0\|_2$ of each method; (b) the mixing parameter β_k of Min-AM and the optimal choice $2/(\mu+L)$; (c) $\|r_k\|_2/\|r_0\|_2$ of AM-I and Min-AM with different β_k .

large. Note that AM-I fails to coincide with Min-AM in the later iterations. It is due to the fact that AM-I needs to solve $(X_k^T R_k)^{-1} X_k^T r_k$ to determine Γ_k , where the matrix inverse operation can have numerical weakness when k is large.

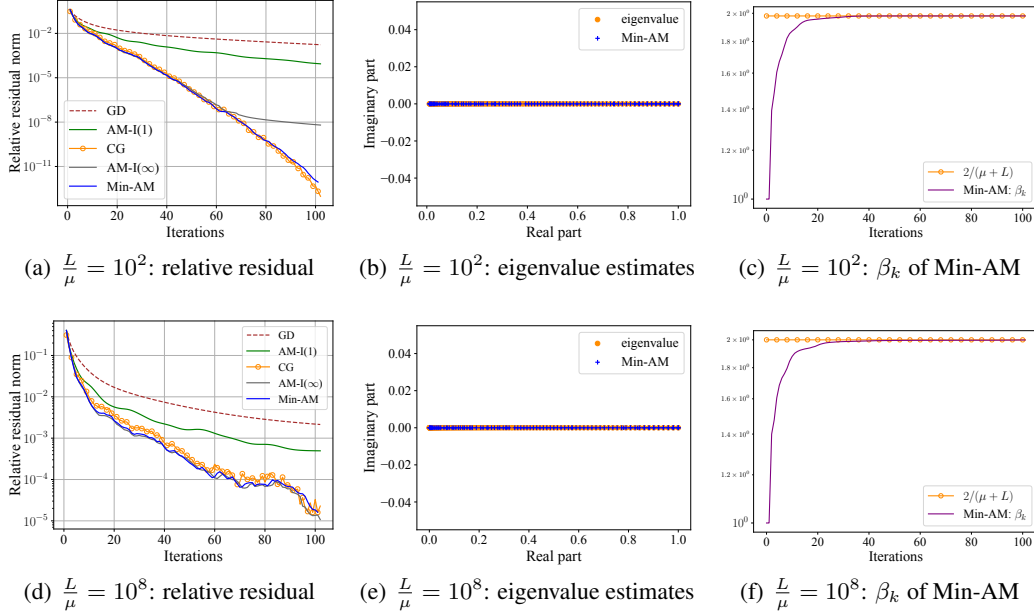


Figure 5: (a)(b)(c): relative residual norm, eigenvalue estimates, and β_k of Min-AM, when $L/\mu = 10^2$. (d)(e)(f): relative residual norm, eigenvalue estimates, and β_k of Min-AM, when $L/\mu = 10^8$. In (b) and (c), “eigenvalue” denotes the exact eigenvalues, and “Min-AM” denotes the eigenvalue estimates computed by Min-AM.

Results about the problem with different condition numbers. In Figure 5, we show the convergence of each method, the eigenvalue estimates from Min-AM, and the β_k of Min-AM, in the tests with different condition numbers characterized by L/μ . The eigenvalues of $A \in \mathbb{R}^{100 \times 100}$ were chosen to be in $(0, 1]$ with equal interval. The results show that Min-AM is competitive with CG and the full-memory AM-I. Min-AM also gives accurate enough estimates of the largest and the smallest eigenvalues, and the β_k is quickly adapted to approximate the optimal value $2/(\mu+L)$.

Cost of the eigenvalue estimation procedure. In our implementation of the eigenvalue estimation procedure in the Min-AM, we used the function “numpy.linalg.eigvals” in NumPy to compute the eigenvalues of T_k constructed by (18), which needs $\mathcal{O}(k^3)$ flops. Hence, the cost of the eigenvalue estimation increases with increasing k , and at the k -th iteration, the ratio of this cost to the total cost is $\mathcal{O}(k^3)/(\mathcal{O}(k^3) + \mathcal{O}(d^2))$, where $\mathcal{O}(d^2)$ is due to the gradient evaluation. To investigate the

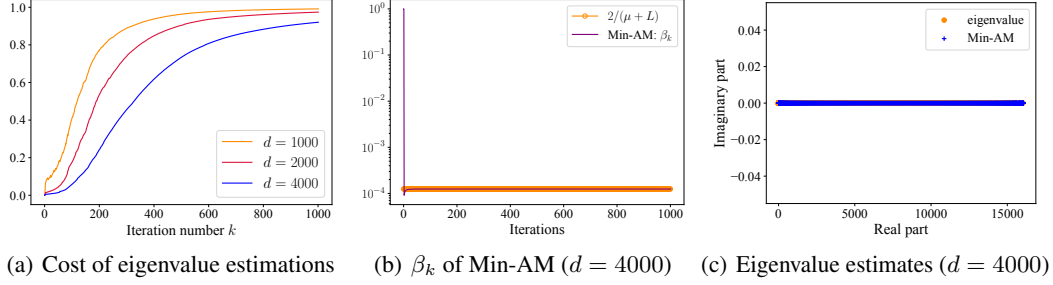


Figure 6: (a) The ratio of the time of eigenvalue estimations to the total running time during iterations; (b) the mixing parameter β_k of Min-AM and the optimal choice $2/(\mu + L)$ when the problem dimension $d = 4000$; (c) exact eigenvalues of the Hessian and eigenvalue estimates computed by Min-AM when the problem dimension $d = 4000$.

practical performance of the eigenvalue estimation, we applied the basic Min-AM to solve problem (172) of dimension $d = 1000, 2000, 4000$, where each $A \in \mathbb{R}^{d \times d}$ was generated following Gaussian distribution, and the maximal iteration number is $max_iter = 1001$ for Min-AM. As can be found in Figure 6(a), the additional computational time incurred by the eigenvalue estimations is marginal when d is large and k is small. Figure 6(b) and Figure 6(c) show the case that $d = 4000$. (The cases $d = 1000$ and $d = 2000$ have similar results.) We find that β_k is quickly adapted to the optimal value $2/(\mu + L)$ within very small number of iterations. So we can adaptively choose β_k in the beginning, and fix the obtained β_k and disable the eigenvalue estimation procedure in the later iterations to reduce the computational cost. Figure 6(c) shows that the eigenvalue estimates computed by Min-AM (at the last iteration) well approximate the exact eigenvalues.

F.2 Regularized logistic regression problem

The regularized logistic regression problem is formulated as

$$\min_{x \in \mathbb{R}^d} f(x) := \frac{1}{T} \sum_{i=1}^T \log(1 + \exp(-y_i x^T \xi_i)) + \frac{w}{2} \|x\|_2^2, \quad (173)$$

where $\xi_i \in \mathbb{R}^d$ is the i -th data sample and $y_i = \pm 1$ is the corresponding label. The regularization parameter $w = 0.1$. We used the datasets “madelon” and “a9a” from LIBSVM [14] that are two-class classifications:

- madelon: training data size: 2000; feature size: 500;
- a9a: training data size: 32561; feature size: 123.

To set proper μ and L for Nesterov’s accelerated gradient (NAG) method and check the eigenvalue estimates from Min-AM, we applied the standard Lanczos algorithm [24] to compute 100 Ritz values of $\nabla^2 f(x^*)$ as the estimates of the true eigenvalues of $\nabla^2 f(x^*)$, where the minimizer x^* was obtained by solving (173). For the test on madelon, Lanczos algorithm gave $\mu = 1.01 \times 10^{-1}$, $L = 1.45$; for the test on a9a, Lanczos algorithm gave $\mu = 1.00 \times 10^{-1}$, $L = 8.64 \times 10^{-1}$.

For the gradient descent (GD), the step size was tuned and set as 1, which was proper for both datasets. We used the Polak-Ribière variant of nonlinear conjugate gradient (NCG), and the step size was determined by line search with cubic interpolation, where the strong Wolfe conditions were checked. The L-BFGS used the Barzilai–Borwein step size as the initial guess of the approximate inverse Hessian. For AM and ST-AM, the mixing parameter was set as 1. The AM-I performed similarly to AM in this test, so we only report the results of AM. For the restarted Min-AM, we set $\tau = 10^{-16}$ for the madelon test and $\tau = 10^{-32}$ for the a9a test. $m = 100$ for both datasets. η was set as a large number since Min-AM converged for these tests. Besides, we computed at most 20 eigenvalue estimates in Min-AM (between two successive restarts) to determine β_k .

Results about the choice of β_k . In Figure 7, we show additional results of the tests on madelon dataset. (1) To test the effect of β_k on Min-AM, we disabled the adaptive choice of β_k and used

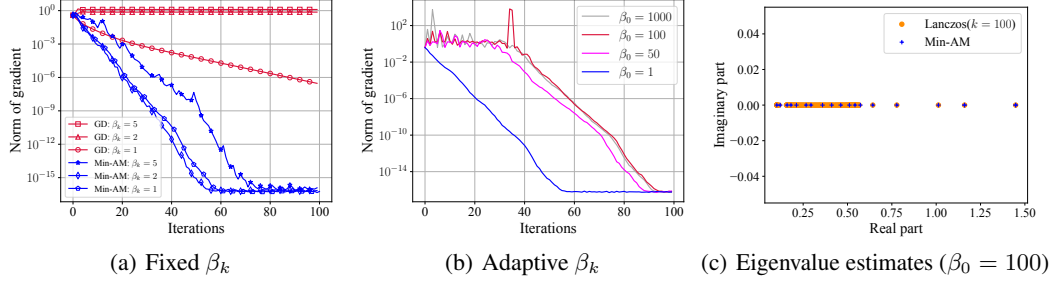


Figure 7: Regularized logistic regression with $w = 0.1$ on madelon dataset. (a) Gradient norms of GD (with constant step size β_k) and Min-AM (with fixed mixing parameter β_k). (b) Min-AM with adaptive choice of β_k ; gradient norms of Min-AM with different initializations of β_k are shown. (c) Min-AM with adaptive choice of β_k ; the Ritz values of $\nabla^2 f(x^*)$ from k -step Lanczos algorithm, and eigenvalue estimates from Min-AM with $\beta_0 = 100$ are shown.

fixed β_k chosen from $\{1, 2, 5\}$. It is observed in Figure 7(a) that GD does not converge if the step size β_k is chosen from $\{2, 5\}$, which suggests that $\theta_k > 1$ in Theorem 2 and Theorem 6. However, the minimization problem on the right-hand side of (112) in Theorem 6 dominates the convergence when m_k is large. It also indicates that Min-AM is less sensitive to β_k than GD. (2) When adaptive choice of β_k is used, we investigate the effect of the initialization β_0 on the convergence and the quality of eigenvalue estimates. It is found in Figure 7(b) that even with the improper initialization $\beta_0 = 50, 100, 1000$, the iterations still finally converge to an acceptable precision. In fact, GD, AM, and Min-AM diverge if using fixed β_k chosen from $\{50, 100, 1000\}$. Hence, Min-AM with adaptive choice of β_k is easy to apply since it discards the requirement of manually tuning the mixing parameter. In Figure 7(c), we find that the eigenvalue estimates can still roughly approximate the largest and the smallest Ritz values computed by the Lanczos algorithm, which verifies Theorem 3.

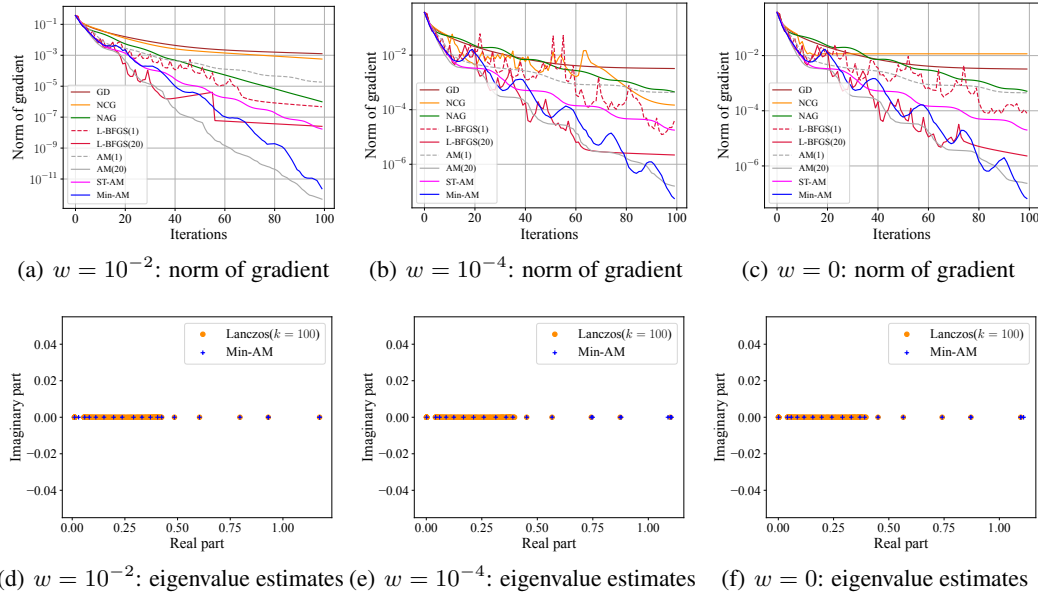


Figure 8: The norm of gradient and eigenvalue estimates for different w . Left column: $w = 10^{-2}$; middle column: $w = 10^{-4}$; right column: $w = 0$. “Lanczos($k = 100$)” denotes the Ritz values computed by Lanczos algorithm, and “Min-AM” denotes the eigenvalue estimates from Min-AM.

Results of the problem with different regularization parameters. We also tested Min-AM for the problem (173) with different settings of the regularization parameter w . The results in Figure 8 also

show that Min-AM significantly improves the convergence of AM(1) and is comparable to AM(20) when w is small. The comparison between the Ritz values computed by Lanczos algorithm and the eigenvalue estimates computed by Min-AM validates Theorem 3. Since the Ritz values approximate the true eigenvalues of $\nabla^2 f(x^*)$, it is expected that Min-AM can give promising eigenvalue estimates of $\nabla^2 f(x^*)$, which accounts for the efficiency of Min-AM.

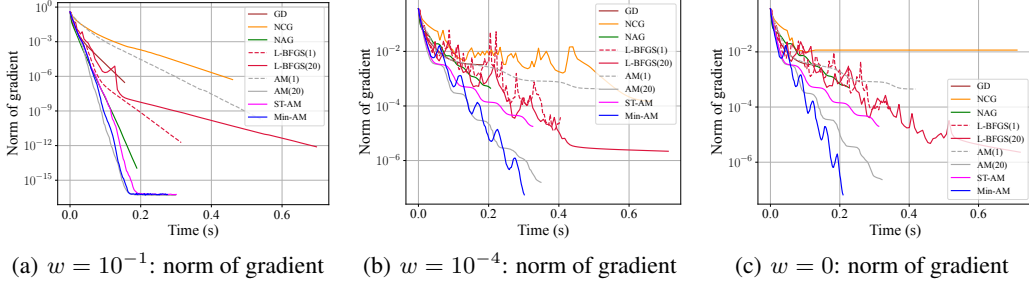


Figure 9: The norm of gradient with respect to running time.

Running time. In Figure 9, we report the convergence of each method with respect to the running time. Min-AM has comparable performance to AM(20) in these tests, while using a smaller memory size.

F.3 Deep neural network training problem

The experiments about the stochastic Min-AM (sMin-AM) focused on training deep neural networks. Since Min-AM can be viewed as a special case of sMin-AM with $\delta_k^{(1)} = \delta_k^{(2)} = 0$, $\alpha_k = 1$, the basic Min-AM was also covered.

F.3.1 Hyperparameter setting of the stochastic Min-AM

The main hyperparameters of sMin-AM are $\alpha_k, \beta_k, \delta_k^{(1)}, \delta_k^{(2)}$. We set $\delta_k^{(1)} = \delta_k^{(2)} = 1$, which ensures that H_k^A is positive definite. We fixed $\alpha_k = 1$ and only tuned $\beta_k \in (0, 1]$ in the experiments.

F.3.2 Experiments on CIFAR-10 and CIFAR-100

The experiments on CIFAR-10 and CIFAR-100 followed the same settings of SAM [62] and ST-AM [63] for direct comparisons. Both CIFAR-10 and CIFAR-100 contain a training dataset (50K images) and a test dataset (10K images), where CIFAR-10 has 10 classes and CIFAR-100 has 100 classes for classification. The basic setting of network training followed the standard setting of training ResNet [30]: The batch size was 128 as commonly suggested; for training with N iterations, the learning rate of the optimizer was decayed at the $(\lfloor \frac{N}{2} \rfloor)$ -th iteration and the $(\lfloor \frac{3}{4}N \rfloor)$ -th iteration. The experiments were run with 3 random seeds and the averaged results along with the standard deviations were reported. The final accuracy on the test dataset was used as the evaluation metric. The accuracy reported in Table 1(a) is the final test accuracy of training with 160 epochs.

We compared sMin-AM with SGDM, Adam, AdaHessian, stochastic AM (SAM), and short-term recurrence AM (ST-AM). SGDM is the default optimizer for training many deep neural networks, such as ResNet [30], WideResNet [66], DenseNet [34], and ResNeXt [64]. Adam is an adaptive learning rate method that uses diagonal approximation to the Hessian based on moving average. AdaHessian also uses adaptive learning rates like Adam, but it relies on Hessian-vector products to obtain the diagonal approximation. SAM and ST-AM are variants of AM. Both can be used to train neural networks. The SAM with memory size m is denoted as SAM(m). ST-AM keeps two vector pairs and has the same memory size as SAM(2).

We tuned the hyperparameters of all the optimizers (including sMin-AM) following the same way for fair comparison. The tuning procedures were conducted on CIFAR-10/ResNet20. For each optimizer, the hyperparameter setting that achieved the highest final test accuracy for training ResNet20 on CIFAR-10 was unchanged in the other tests of training neural networks on CIFAR.

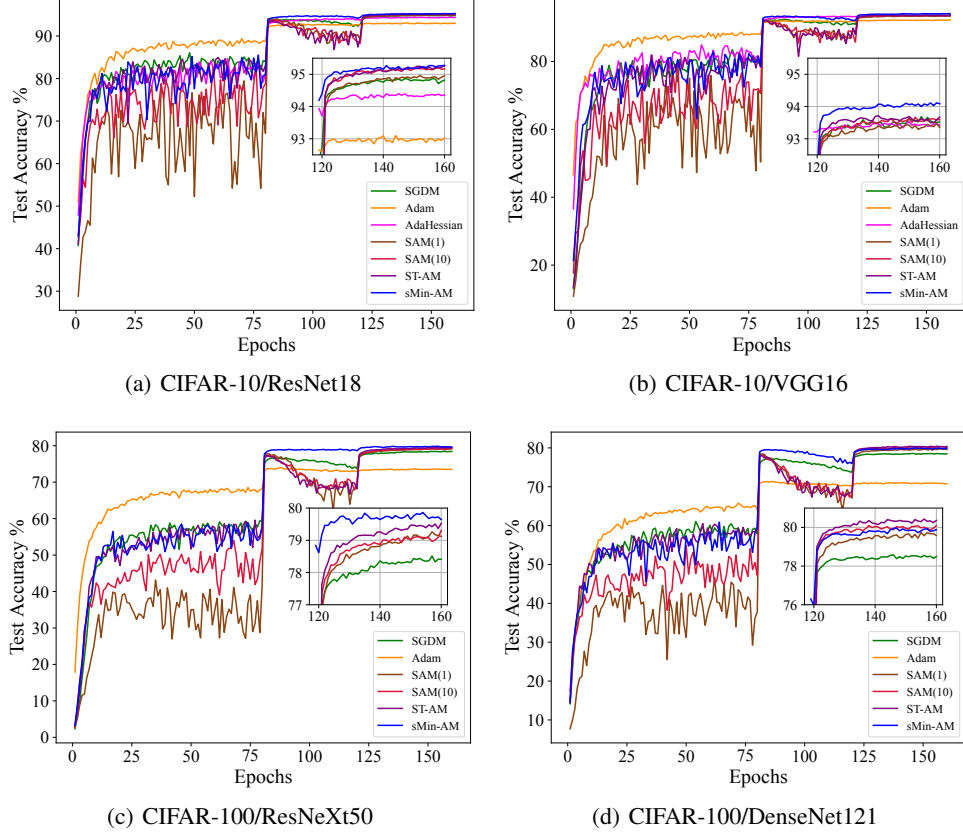


Figure 10: Test accuracy during the network training on CIFAR-10 and CIFAR-100.

For SGDM, the momentum was set as 0.9, which is the default setting in the literature [30, 34]. The initial learning rate and weight decay were tuned and were set as 0.1, 0.0005, respectively. The learning rate decay was 0.1. This setting is also recommended in WideResNet [66].

For Adam, the initial learning rate and weight decay were tuned and set as 0.001, 0.0005, respectively. The momentum terms were $\beta_1 = 0.9$ and $\beta_2 = 0.999$ as commonly suggested [65, 68]. The learning rate decay was 0.1.

For AdaHessian, the initial learning rate was 0.15, and the momentum terms $\beta_1 = 0.9$, $\beta_2 = 0.999$; $\text{eps} = 1 \times 10^{-4}$, and the hessian-power was 1. The weight decay was 0.0005/0.15, and the learning rate decay was 0.1.

For SAM, the initial mixing parameter $\beta_0 = 1$, the initial damping term $\alpha_0 = 1$, and the regularization parameter $c_1 = 0.01$. The weight decay was 0.0015. The decay rate for β_k, α_k was 0.06.

For ST-AM, the initial mixing parameter $\beta_0 = 1$, the initial damping term $\alpha_0 = 1$, and the regularization parameters were $c_1 = 1$, $c_2 = 1 \times 10^{-7}$. The weight decay was 0.001, and the decay rate for β_k, α_k was 0.1.

For sMin-AM, the initial mixing parameter $\beta_0 = 0.2$, the damping term $\alpha_k = 1$, and the regularization parameters were $\delta_k^{(1)} = 1$, $\delta_k^{(2)} = 1$. The weight decay was 0.0015, and the decay rate for β_k was 0.1.

Note that our hyperparameter settings of the baseline methods were the same as those in SAM [62] and ST-AM [63], so we used their results for reference.

Convergence behaviour in the training process. In Figure 10, we plot the test accuracy of training four networks on CIFAR-10 and CIFAR-100, for each optimizer. It shows Adam has fast convergence in the beginning, but stagnates in the later training. Compared with SGDM, SAM, and ST-AM, it is

observed that sMin-AM converges faster to an acceptable accuracy. The process of sMin-AM is also much more stable than SAM(1) and SAM(10).

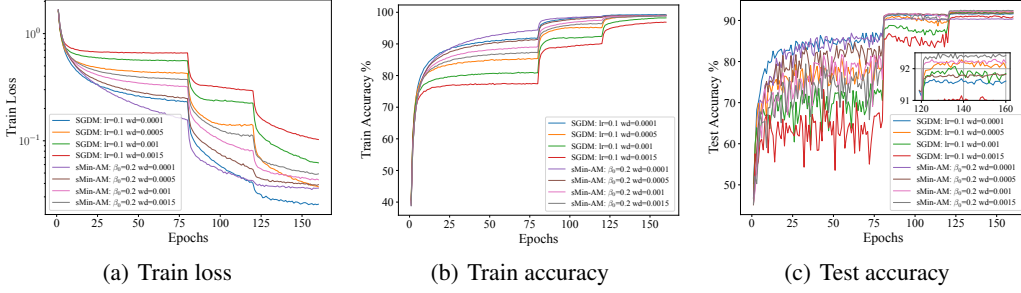


Figure 11: Effect of weight decay on training CIFAR-10/ResNet20. “lr” and “wd” are abbreviations of learning rate and weight decay.

Effect of weight decay. Figure 11 shows the effect of weight decay on the training process. It suggests that a larger weight decay can slow down the training in terms of training loss. When using the same weight decay, sMin-AM often has faster convergence than SGDM. The results also justify our settings of the weight decay.

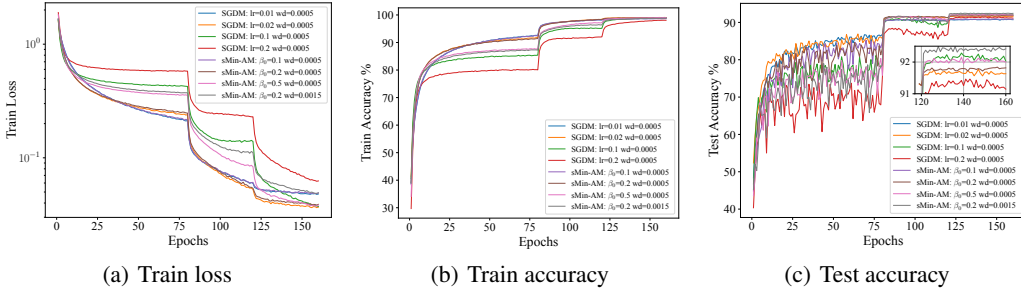


Figure 12: Effect of learning rate or mixing parameter on training CIFAR-10/ResNet20. “lr” and “wd” are abbreviations of learning rate and weight decay.

Effect of learning rate/mixing parameter. Figure 12 shows the behaviours of SGDM with different learning rate settings and sMin-AM with different mixing parameters. Small learning rate or small mixing parameter can lead to fast convergence in the beginning, but the training may stagnate early and the test performance is not satisfactory.

Effects of $\delta_k^{(1)}$, $\delta_k^{(2)}$, and α_k . In Figure 13, we show the results of sMin-AM with different settings of other hyperparameters, i.e., $\delta_k^{(1)}$, $\delta_k^{(2)}$, and α_k . It is observed that sMin-AM is more sensitive to $\delta_k^{(2)}$ than $\delta_k^{(1)}$ and α_k , which may be due to the fact that $\delta_k^{(2)}$ directly affects the regularization in the Min-AM update (Line 10 in Algorithm 4). When choosing $\delta_k^{(2)} \geq 1$, the effect of $\delta_k^{(2)}$ is also minor.

As a result, the most critical hyperparameters that affect the effectiveness of sMin-AM are the mixing parameter β_k , the regularization parameter $\delta_k^{(2)}$, and weight decay. The weight decay is a common hyperparameter that needs to be tuned for each optimizer; for $\delta_k^{(2)}$, we can choose $\delta_k^{(2)} \geq 1$ to ensure $H_k^A \succeq 0$. So except for the weight decay, we only tuned the mixing parameter in the experiments.

Check of the choices of $\delta_k^{(1)}$ and $\delta_k^{(2)}$. From Lemma 2, we know $\|H_k^A - \beta_k I\|_2 \leq -\rho_k \|p_k\|_2^2 - 2\beta_k \rho_k \|p_k\|_2 \|q_k\|_2 + \beta_k \rho_k^2 \|p_k\|_2^2 \|q_k\|_2^2$ assuming $\rho_k \leq 0$. In Figure 14, we plot the value of this bound. It is found that $\|H_k^A - \beta_k I\|_2 \leq \beta_k C_2$ for some constant C_2 , which justifies our choice of the regularization parameters.

Train neural networks with fewer epochs and less time. We also tested sMin-AM for training different epochs. Figure 15 shows the comparison between sMin-AM and SGDM for training 80,

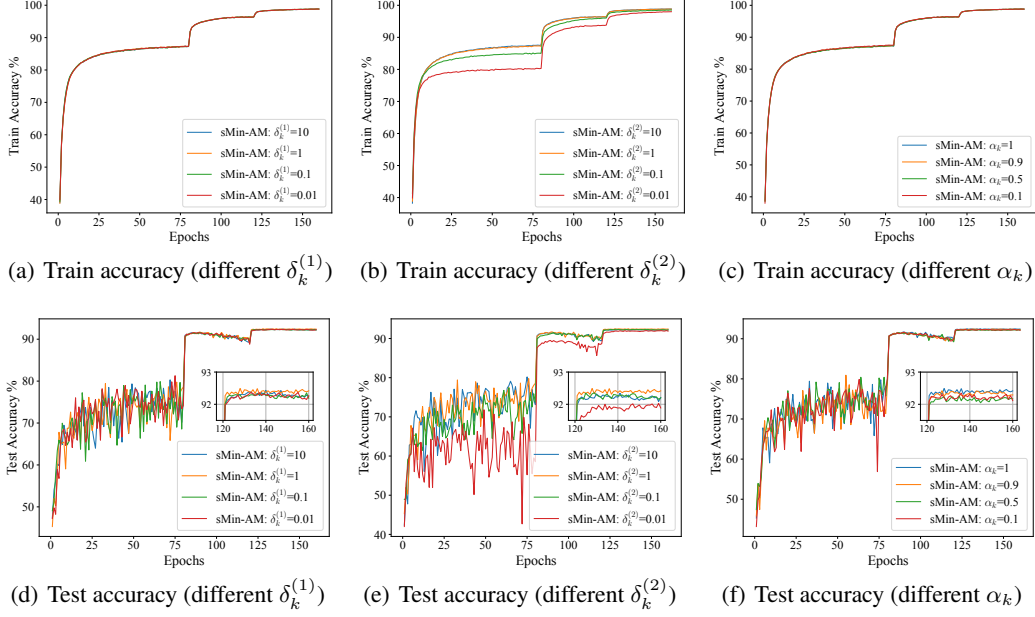


Figure 13: Effects of $\delta_k^{(1)}$, $\delta_k^{(2)}$, and α_k on sMin-AM for training CIFAR-10/ResNet20. The default setting is $\delta_k^{(1)} = 1$, $\delta_k^{(2)} = 1$, $\alpha_k = 1$, $\beta_0 = 1$, and weight decay is 0.0015. When one hyperparameter was inspected, the other hyperparameters were set as default.

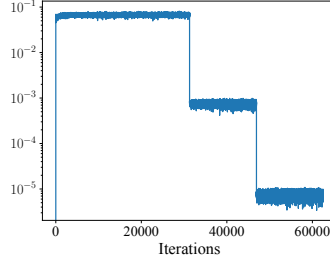


Figure 14: The value of $-\rho_k \|p_k\|_2^2 - 2\beta_k \rho_k \|p_k\|_2 \|q_k\|_2 + \beta_k \rho_k^2 \|p_k\|_2^2 \|q_k\|_2^2$ during training CIFAR-10/ResNet20.

120, and 160 epochs. The results show that sMin-AM can attain comparable accuracy to SGDM in fewer epochs, which infers the faster training process of sMin-AM. In Table 2, by setting the final test accuracy (160 epochs) of SGDM as the baseline, we report the memory cost, per-epoch time, total training epochs, total training time of SAM(10), ST-AM, and sMin-AM to achieve a comparable accuracy to SGDM (within 0.05% difference). The attained final accuracy is also shown. It is found that sMin-AM largely reduces the memory cost of SAM(10) and can achieve comparable results to those of SGDM using less training time.

Discussion about the computational cost and memory cost. The per-epoch computational cost and memory cost are closely related to three factors: the network architecture, the training data, and the optimizer. From Table 2, the effects of the network architecture and the optimizer are clear. For the effect of training data, we consider the batch size. Given a specific network, the cost is composed of two parts: (i) updating network parameters by the optimizer; (ii) other necessary computations, such as the forward and back propagations of the neural network, data transfers between memory and disks, etc., where additional memory and processing time are required. If the cost of Part (ii) only occupies a small proportion of the total cost, the cost incurred by the optimizer will be of great importance. In Figure 16, we plot the ratios of memory/per-epoch time of SAM(10)/ST-AM/sMin-AM to that of SGDM, for training ResNeXt50 on CIFAR-100 with different batch sizes (16, 32, 64, and 128).

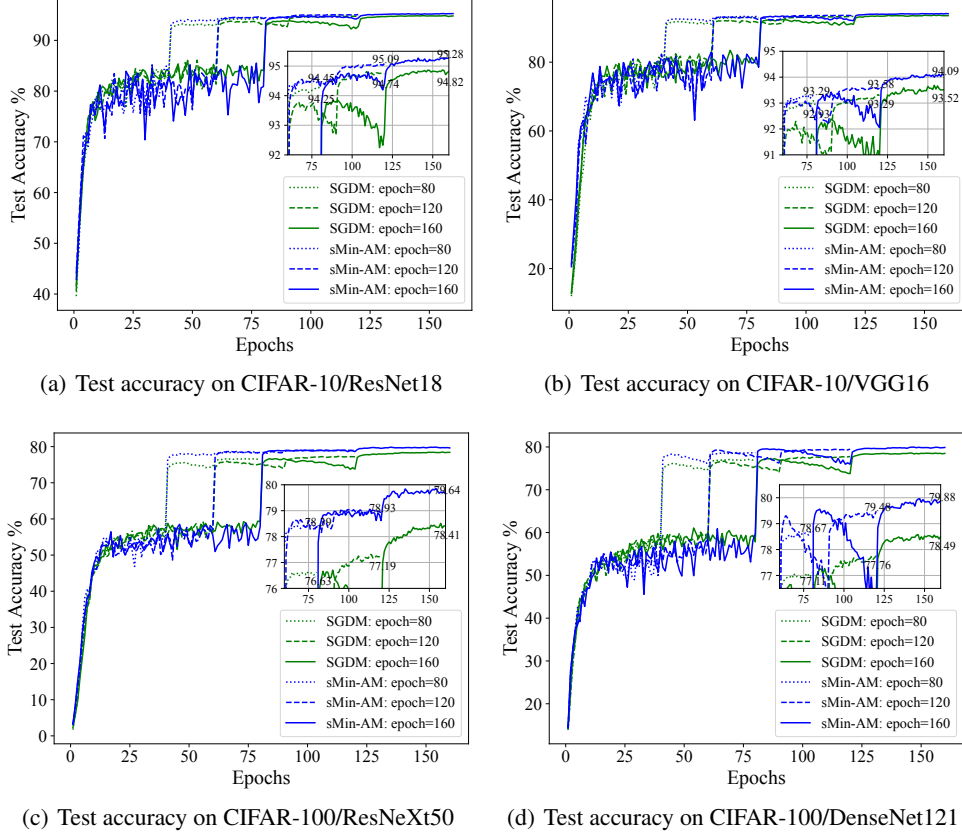


Figure 15: Comparison between SGDM and sMin-AM on training deep neural networks for 80, 120, and 160 epochs.

Table 2: The cost and final test accuracy compared with SGDM. Notations “m”, “t/e”, “e”, “t”, “a” are abbreviations of memory, per-epoch time, training epochs, total training time, accuracy.

Cost (\times SGDM) & accuracy	CIFAR-10/ResNet18					CIFAR-10/VGG16				
	m	t/e	e	t	a(%)	m	t/e	e	t	a(%)
SGDM	1.00	1.00	1.00	1.00	94.82	1.00	1.00	1.00	1.00	93.52
SAM(10)	1.73	1.78	0.56	1.00	94.81	2.51	2.59	1.00	2.59	93.59
ST-AM	1.05	1.46	0.56	0.82	94.84	1.55	1.91	0.88	1.67	93.56
sMin-AM	1.01	1.15	0.56	0.64	94.85	1.35	1.25	0.63	0.78	93.56

Cost (\times SGDM) & accuracy	CIFAR-100/ResNeXt50					CIFAR-100/DenseNet121				
	m	t/e	e	t	a(%)	m	t/e	e	t	a(%)
SGDM	1.00	1.00	1.00	1.00	78.41	1.00	1.00	1.00	1.00	78.49
SAM(10)	1.30	1.16	0.50	0.58	78.37	1.16	1.19	0.50	0.60	78.84
ST-AM	1.04	1.07	0.50	0.54	78.39	1.01	1.11	0.50	0.55	78.90
sMin-AM	1.03	1.00	0.50	0.50	78.39	1.01	1.09	0.50	0.55	78.67

When a smaller batch size is used, the Part (ii) cost reduces and the proportion of Part (i) cost in the total cost increases, thus the performance improvement of sMin-AM over SAM(10)/ST-AM is more significant. To investigate the effect of model size on the cost, we tested SAM(10)/ST-AM/sMin-AM on CIFAR-10/ResNet18 and CIFAR-10/ResNet50, where the batch size is 16 so that the Part (ii) cost is not large. In Table 3, it is observed that for a network of larger scale (ResNet50), sMin-AM is more advantageous than SAM(10)/ST-AM in terms of memory cost.

The difference between CIFAR-10 and CIFAR-100 has only a minor effect on the cost. For a given network architecture, e.g., VGG16, its implementations for CIFAR-10 and CIFAR-100 differ in the

dimension of the last linear layer. Since the parameters of this linear layer usually occupy a very small portion of the whole parameters, the costs (memory and per-epoch running time) on CIFAR-10 and CIFAR-100 are roughly the same, as shown in Table 4.

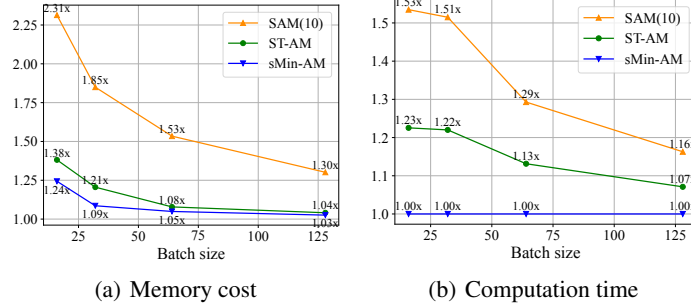


Figure 16: Memory and per-epoch time for training CIFAR-100/ResNeXt50 with different batch sizes, where those of SGDM are set as the units.

Table 3: The memory and per-epoch time of SAM(10)/ST-AM/sMin-AM, where the batch size is 16. The memory and per-epoch time of SGDM are set as the units; memory and per-epoch time are abbreviated as “m” and “t/e”, respectively.

Cost (\times SGDM)	CIFAR-10/ResNet18		CIFAR-10/ResNet50	
	m	t/e	m	t/e
SAM(10)	2.44	2.66	3.01	2.63
ST-AM	1.59	1.91	1.79	1.86
sMin-AM	1.41	1.21	1.52	1.15

Table 4: The memory and per-epoch time of SAM(10)/ST-AM/sMin-AM, where the batch size is 128. The memory and per-epoch time of SGDM are set as the units; the first element and the second element in (\cdot, \cdot) denote the memory and per-epoch time, respectively.

Cost (\times SGDM)	CIFAR-10				CIFAR-100			
	ResNet18	VGG16	ResNeXt50	DenseNet121	ResNet18	VGG16	ResNeXt50	DenseNet121
SAM(10)	(1.73, 1.78)	(2.51, 2.59)	(1.30, 1.18)	(1.16, 1.21)	(1.84, 1.85)	(2.52, 2.53)	(1.30, 1.16)	(1.16, 1.19)
ST-AM	(1.05, 1.46)	(1.55, 1.91)	(1.04, 1.06)	(1.02, 1.08)	(1.12, 1.47)	(1.56, 1.84)	(1.04, 1.07)	(1.01, 1.11)
sMin-AM	(1.01, 1.15)	(1.35, 1.25)	(1.02, 1.00)	(1.01, 1.01)	(1.08, 1.09)	(1.36, 1.19)	(1.03, 1.00)	(1.01, 1.09)

F.3.3 Experiment of training ResNet50 on ImageNet

We applied sMin-AM to train ResNet50 on the ImageNet dataset, which contains 1.2M images for training and 50K images for test. The code was based on the example from PyTorch³. We followed the standard process of training ResNet [30]. Four GPUs were used to conduct the experiment. SGDM was used for comparison. The batch size was 256. The weight decay for each method was 0.0001. We trained ResNet50 for 90 epochs. The learning rate of SGDM and the mixing parameter of sMin-AM were decayed by 0.1 at the 30th and 60th epochs. For SGDM, the momentum was 0.9 and the learning rate was 0.1, which was the standard setting [30, 66, 34, 64]. For sMin-AM, we set the initial mixing parameter $\beta_0 = 0.5$, and the other hyperparameters were kept unchanged.

Comparisons with SAM and ST-AM. We conducted tests to compare the accuracy of sMin-AM with that of SAM/ST-AM, and results in Figure 17 show that sMin-AM also improves the test accuracy of SAM(10)/ST-AM.

F.3.4 An additional experiment: adversarial training

We conduct an additional experiment to further compare the effectiveness of SGDM, ST-AM, and sMin-AM for training deep neural networks. The considered problem is *adversarial training*, which

³ <https://github.com/pytorch/examples/tree/master/imagenet>.

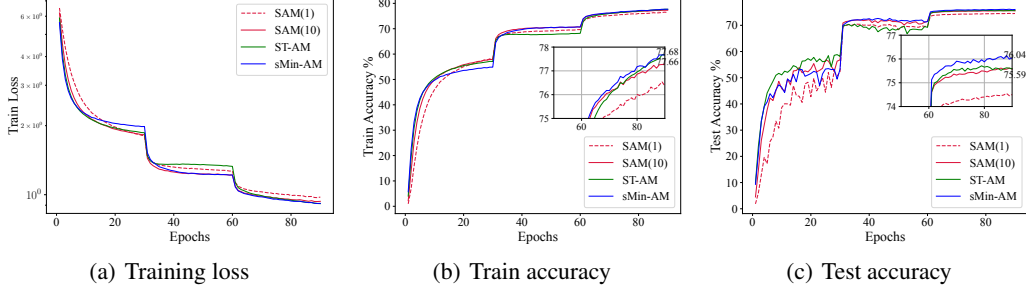


Figure 17: Comparisons between SAM, ST-AM, and sMin-AM for training ResNet50 on ImageNet. The results of the final accuracy of ST-AM and sMin-AM are shown in (b) and (c).

tries to solve the following minimax problem

$$\min_{x \in \mathbb{R}^d} \frac{1}{T} \sum_{i=1}^T \max_{\|\bar{\xi}_i - \xi_i\|_2 \leq \epsilon} f_{\bar{\xi}_i}(x), \quad (174)$$

where $\bar{\xi}_i$ is the adversarial data sample in the ϵ -ball centered at the data sample ξ_i . We used the standard PGD adversarial training process [40]: using projection gradient descent to solve the maximization problem in (174), and applying SGDM/ST-AM/sMin-AM to solve the minimization problem in (174). The tests were conducted on CIFAR-10/WideResNet34-10 and CIFAR-100/DenseNet121, where for the training dataset of CIFAR-10/CIFAR-100, 5K images were randomly selected as the validation dataset and the other 45K images were used as the new training dataset. We trained the neural network for 200 epochs, where the learning rate/mixing parameter was decayed at the 100th and the 150th epochs. The best checkpoint model on the validation dataset was chosen to be evaluated on the test dataset. For SGDM and ST-AM, we used the recommended setting of hyperparameters in [26, 63]. For sMin-AM, we set $\beta_0 = 0.3$, and other hyperparameters were kept unchanged.

Table 5: Clean test accuracy (%) and robust test accuracy (%) on adversarial training.

Optimizer	CIFAR-10/WideResNet34-10				CIFAR-100/DenseNet121			
	Clean	FGSM	PGD-20	C&W $_{\infty}$	Clean	FGSM	PGD-20	C&W $_{\infty}$
SGDM	85.48	66.42	54.00	53.28	59.45	39.75	30.91	29.02
ST-AM	85.79	66.43	53.46	52.88	60.48	40.39	31.19	29.56
sMin-AM	85.76	67.49	54.67	54.05	60.04	40.50	31.71	29.95

The training process of each optimizer is shown in Figure 18, where we plot the accuracy on validation dataset (called as validation accuracy), and the accuracy on adversarial data samples generated by the PGD-10 attack [40] on validation dataset (called as robust validation accuracy). It is observed that the best robust validation accuracy of sMin-AM is higher than that of SGDM/ST-AM, which indicates that better checkpoint model can be obtained by sMin-AM.

In Table 5, we report two types of test accuracy: clean test accuracy, where the clean test data was used for classification; robust test accuracy, where corrupted test data was used for classification. Three attacking methods were used for the robust test accuracy evaluation: FGSM [25], PGD-20 [40], and C&W $_{\infty}$ attack [13]. It shows that sMin-AM significantly improves SGDM, and outperforms ST-AM in terms of robust test accuracy, which is a desirable property in adversarial training.

G Limitations

We focus on developing the variant of AM with minimal memory size for solving optimization problems. The properties of Min-AM for deterministic optimization are established in the smooth case. The non-smooth optimization is important and one possible direction is to adapt Min-AM to the framework of proximal quasi-Newton methods [15, 7]. We leave it as an extension of Min-AM in the future work.

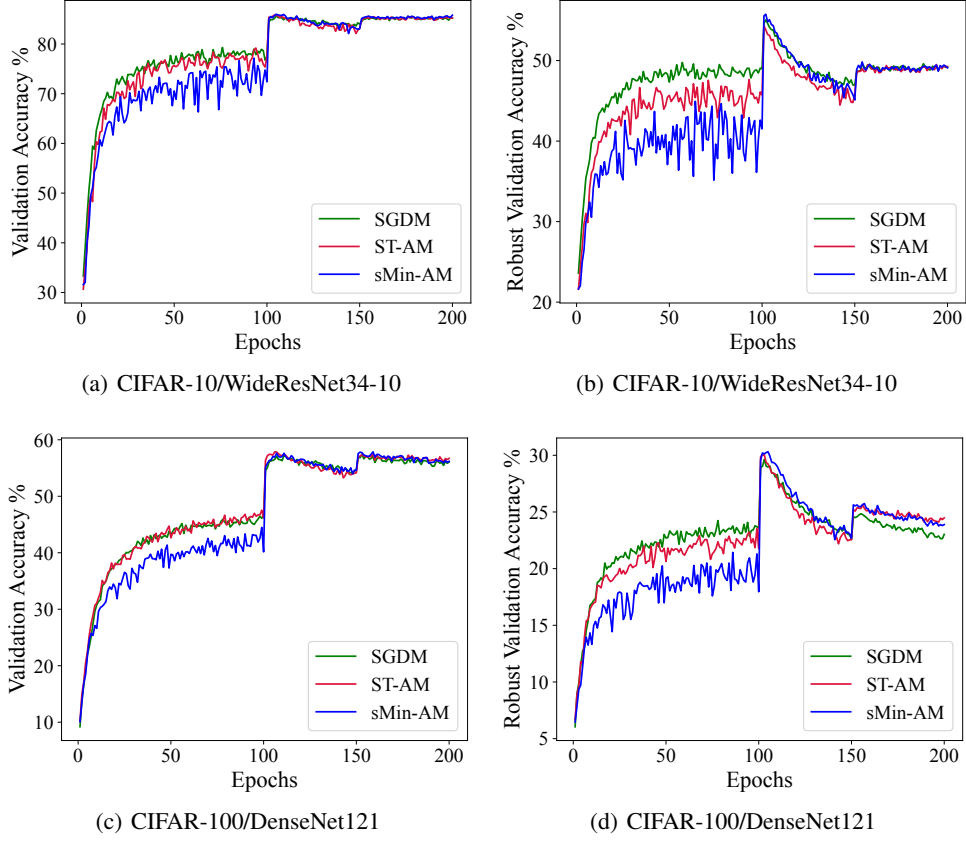


Figure 18: Left: accuracy on the clean validation dataset (validation accuracy). Right: accuracy on the validation dataset attacked by PGD-10 (robust validation accuracy).

Min-AM exploits the symmetry of Hessian in solving optimization problems. For general fixed-point problems, the Jacobian can be non-symmetric. However, in this case, the short-term recurrence algorithms that are equivalent to the full-memory methods generally do not exist even in solving linear systems [53].