

# 关于 Anderson 混合的研究进展\*

包承龙<sup>1</sup>, 韦福超<sup>2</sup>

- 清华大学丘成桐数学科学中心, 北京 100084
- 清华大学计算机科学与技术系, 北京 100084

**摘要:** Anderson 混合是一种经典的外推方法, 它能利用历史迭代信息加速定点迭代的收敛, 在科学计算和机器学习得到了成功的应用. 由于 Anderson 混合在实践中经常表现出优越的数值性能, 在各类应用中围绕 Anderson 混合的算法设计和理论分析成为近几年的研究热点. 本文综述关于 Anderson 混合的研究进展, 重点介绍基于 Anderson 混合的新算法.

**关键词:** Anderson 混合; 定点迭代; Krylov 子空间方法; 拟 Newton 法

**中图分类号:** O221.2 **文献标志码:** A **文章编号:** 2097-0137(2023)05-0059-08

## Research advance on Anderson mixing

BAO Chenglong<sup>1</sup>, WEI Fuchao<sup>2</sup>

- Yau Mathematical Sciences Center, Tsinghua University, Beijing 100084, China
- Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China

**Abstract:** Anderson mixing is a classical extrapolation method. It can make use of the information in historical iterations to accelerate the convergence of fixed-point iterations, and has been successfully applied in scientific computing and machine learning. Since Anderson mixing often exhibits superior numerical performance in practice, the algorithm design and theoretical analysis around Anderson mixing in various applications have become hot topics in recent years. This article reviews the research advance on Anderson mixing, and highlights new algorithms based on Anderson mixing.

**Key words:** Anderson mixing; fixed-point iteration; Krylov subspace method; quasi-Newton method

由 Anderson(1965)提出的 Anderson 混合(AM, Anderson mixing)最早被用于非线性积分方程的计算, 现已成为加速定点迭代的一种经典算法. 在量子化学领域, AM 又被称为 Pulay 混合(Pulay, 1980)或 DIIS 方法(Rohwedder et al., 2011), 在自洽场迭代的加速中发挥了重要作用(Arora et al., 2017). AM 本质是一种外推方法(Brezinski et al., 2018; Anderson, 2019), 它通过对历史迭代外推, 生成与定点迭代不同的新迭代序列. 由于 AM 通常能够显著减少迭代过程收敛到定点的迭代次数, 因此和定点迭代相比, 当定点算子的计算开销很大时, AM 算法能够节省大量的计算时间. 所以, AM 在科学计算中也常被称为 Anderson 加速(Walker et al., 2011). 近几年来, 得益于其实现的简易性和优异的数值表现, AM 在科学计算和机器学习领域得到了广泛的关注, 研究者们成功地将 AM 应用于各种定点问题的求解当中, 例如, 解 Navier-Stokes 方程(Pollock et al., 2019)、解地震波反演问题(Yang, 2021)、加速机器学习中的 EM 算法(Henderson et al., 2019)、强化学习训练(Sun et al., 2021)等. 此外, AM 的理论性质也引起计算数学界的极大兴趣, 对 AM 在定点问题中的收敛性给出理论分析仍是目前一个重要的研究问题(Toth et al., 2015; Evans et al., 2020;

\* 收稿日期: 2023-05-11 录用日期: 2023-06-21 网络首发日期: 2023-08-30

基金项目: 国家重点研发计划(2021YFA1001300); 国家自然科学基金(12271291)

作者简介: 包承龙(1989年生), 男; 研究方向: 图像处理模型与优化算法; E-mail: clbao@mail.tsinghua.edu.cn

Bian et al., 2021).

先简要介绍 AM 的基本迭代格式. 考虑定点问题

$$\mathbf{x} = g(\mathbf{x}), \quad (1)$$

其中  $\mathbf{x} \in \mathbb{R}^d$ ,  $g: \mathbb{R}^d \rightarrow \mathbb{R}^d$ . 如果  $g$  是收缩算子, 那么由压缩映射原理, 定点迭代

$$\mathbf{x}_{k+1} = g(\mathbf{x}_k), \quad k = 0, 1, \dots, \quad (2)$$

收敛. 为加速迭代(2), AM 对历史迭代序列进行外推来生成新的迭代值. 具体地, 设第  $k$  次迭代用到的历史序列的长度为  $m_k$ , AM 使用下式更新得到

$$\mathbf{x}_{k+1} = (1 - \beta_k) \sum_{j=0}^{m_k} \alpha_k^{(j)} \mathbf{x}_{k-m_k+j} + \beta_k \sum_{j=0}^{m_k} \alpha_k^{(j)} g(\mathbf{x}_{k-m_k+j}), \quad (3)$$

其中  $\beta_k$  为混合参数;  $\{\alpha_k^{(j)}\}_{j=0}^{m_k}$  为外推系数, 由求解以下带约束的最小二乘问题得到:

$$\min_{\{\alpha_k^{(j)}\}_{j=0}^{m_k}} \left\| \sum_{j=0}^{m_k} \alpha_k^{(j)} (g(\mathbf{x}_{k-m_k+j}) - \mathbf{x}_{k-m_k+j}) \right\|_2, \quad \text{s.t.} \quad \sum_{j=0}^{m_k} \alpha_k^{(j)} = 1. \quad (4)$$

随后将看到问题(4)实际是一个残差极小化问题, 能够让外推系数的确定符合某个最优条件. 如果限制外推系数均非负, 即  $\alpha_k^{(j)} \geq 0$ ,  $j = 0, \dots, m_k$ , 就得到 EDIIS 方法(Kudin et al., 2002).

相较于定点迭代, AM 的迭代更新过程的主要开销在于存储历史序列并对之完成外推计算, 而并不需要多余的关于  $g$  的计算. 在之后的研究中, 人们基于 AM 的迭代格式发展得到更多的算法, 这些算法在一些更具体的问题求解中比原始的 AM 有更好的性质和表现.

由于定点问题广泛存在于科学与工程各个领域, AM 有相当广阔的适用场景, 因此在各类应用中围绕 AM 的算法设计和理论分析是目前科学界的前沿热点, 本文将对关于 AM 的研究进展加以介绍. 接下来, 本文深入剖析 AM 的迭代格式, 随后重点介绍关于 AM 的几个改进算法, 算法包括正则化的 AM、随机 AM、短递归 AM 和具有极小内存开销的 AM 算法, 这些算法能够在很大程度上拓展 AM 的应用范围.

这里给出文中的符号定义: 符号  $\Delta$  为前向差分符号, 例如  $\Delta \mathbf{x}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ ; 符号  $\dagger$  表示取 Penrose-Moore 逆. 对任意矩阵  $A$ ,  $\text{range}(A)$  表示由  $A$  的列向量张成的子空间. 矩阵范数  $\|\cdot\|_{F(W)}$  的定义为对任意  $X \in \mathbb{R}^{d \times d}$ , 有  $\|X\|_{F(W)} = \|W^{1/2} X W^{1/2}\|_F$ .

## 1 基础算法

本节在投影-混合框架下分析 AM 的迭代格式, 给出第一类 AM 方法, 并介绍已有的结论.

对于求解问题(4), 一种方式是通过 Lagrange 乘子法求解, 另一种方式是将其转换为无约束问题. 具体地, 定义  $\mathbf{x}_k$  处的残差为  $\mathbf{r}_k = g(\mathbf{x}_k) - \mathbf{x}_k$ , 历史序列被存储为两个矩阵  $X_k, R_k \in \mathbb{R}^{d \times m_k}$  ( $m_k \geq 1$ ):

$$X_k = (\Delta \mathbf{x}_{k-m_k}, \Delta \mathbf{x}_{k-m_k+1}, \dots, \Delta \mathbf{x}_{k-1}), \quad R_k = (\Delta \mathbf{r}_{k-m_k}, \Delta \mathbf{r}_{k-m_k+1}, \dots, \Delta \mathbf{r}_{k-1}). \quad (5)$$

AM 的迭代格式可以被分解为投影步和混合步:

$$\bar{\mathbf{x}}_k = \mathbf{x}_k - X_k \Gamma_k \quad (\text{投影步}), \quad \bar{\mathbf{r}}_k = \mathbf{r}_k - R_k \Gamma_k, \quad \mathbf{x}_{k+1} = \bar{\mathbf{x}}_k + \beta_k \bar{\mathbf{r}}_k \quad (\text{混合步}), \quad (6)$$

其中  $\bar{\mathbf{x}}_k$  和  $\bar{\mathbf{r}}_k$  分别为中间步和中间残差, 二者通过混合步导出  $\mathbf{x}_{k+1}$ .  $\Gamma_k$  由求解以下的最小二乘问题确定:

$$\Gamma_k = \arg \min_{\Gamma \in \mathbb{R}^{m_k}} \|\mathbf{r}_k - R_k \Gamma\|_2. \quad (7)$$

因为  $\bar{\mathbf{r}}_k = \mathbf{r}_k - R_k \Gamma_k$ , 所以问题(7)可以被视作一个残差极小化问题. 令  $\Gamma_k = (\Gamma_k^{(1)}, \dots, \Gamma_k^{(m_k)})^T \in \mathbb{R}^{m_k}$ , 那么外推系数  $\{\alpha_k^{(j)}\}_{j=0}^{m_k}$  能由  $\Gamma_k$  得到:  $\alpha_k^{(0)} = \Gamma_k^{(1)}$ ,  $\alpha_k^{(j)} = \Gamma_k^{(j+1)} - \Gamma_k^{(j)}$  ( $j = 1, \dots, m_k - 1$ ),  $\alpha_k^{(m_k)} = 1 - \Gamma_k^{(m_k)}$ . 因为有  $\mathbf{r}_k - R_k \Gamma_k = \sum_{j=0}^{m_k} \alpha_k^{(j)} \mathbf{r}_{k-m_k+j}$ , 所以通过计算可以验证由式(6)~(7)给出的迭代与由式(3)~(4)给出的迭代是完全相同的.

注意到求解问题(7)等价于要求  $\bar{\mathbf{r}}_k = \mathbf{r}_k - R_k \Gamma_k \perp \text{range}(R_k)$ , 因此  $\Gamma_k$  的计算是一个投影过程. 如果使用投影条件  $\bar{\mathbf{r}}_k = \mathbf{r}_k - R_k \Gamma_k \perp \text{range}(X_k)$ , 就导出第一类 AM 方法(Fang et al., 2009). 与之区分, 原始的 AM 方法被称为第二类 AM 方法. 两类方法的迭代格式可以被写成

$$\mathbf{x}_{k+1} = G_k \mathbf{r}_k = \mathbf{x}_k + \beta_k \mathbf{r}_k - (X_k + \beta_k R_k) \Gamma_k,$$

其中  $G_k, \Gamma_k$  由各自的投影条件确定. Fang et al. (2009) 指出两类 AM 方法实际都为多重割线的拟 Newton 法. 具体而言, 第一类 AM 的  $G_k$  满足  $G_k = J_k^{-1}$ , 其中  $J_k$  求解了  $\min_j \|J - \beta_k^{-1} I\|_F$  s.t.  $JX_k = -R_k$ ; 第二类 AM 的  $G_k$  求解了  $\min_c \|G - \beta_k I\|_F$  s.t.  $GR_k = -X_k$ .

在 AM 的使用中需要确定历史序列长度  $m_k$ . 一种做法是取  $m_k = k$ , 即使用全部的历史信息, 因此这被称为全记忆的方法; 另一种做法是取  $m_k = \min\{m, k\}$ , 即使用最近  $m$  步的迭代信息, 从而限制内存占用, 这被称为有限内存的方法, 将第一类 AM 和第二类 AM 分别记为 AM-I( $m$ ) 和 AM-II( $m$ ). 由于 AM 的外推计算的开销在  $m_k$  较大时较为可观, 因此一种节省开销并保留一定的 AM 加速效果的方式是使用定点迭代和 AM 交替迭代 (Pratapa et al., 2016; Suryanarayana et al., 2019). 在该方案中, 每  $p$  步迭代中先执行  $(p-1)$  步定点迭代, 再执行 1 步 AM 迭代.

关于 AM 的理论分析主要分为两部分, 一个是线性定点问题中全记忆的 AM 和 Krylov 子空间方法的关系, 另一个是非线性定点问题中有限内存的第二类 AM 的收敛性. 以下加以叙述.

对于求解线性方程组, 两类 AM 与 Arnoldi 方法 (Saad, 1981)、GMRES 方法 (Saad et al., 1986) 这两类典型的 Krylov 子空间方法有着本质的联系. 设问题 (1) 中  $g(x) = (I - A)x + b$ ,  $A \in \mathbb{R}^{d \times d}$  非奇异,  $b \in \mathbb{R}^d$ . 令  $\{x_k^I\}$  和  $\{x_k^{II}\}$  分别为全记忆的第一类和全记忆的第二类 AM 生成的序列, 令  $\{x_k^A\}$  和  $\{x_k^G\}$  为 Arnoldi 方法和 GMRES 生成的序列. 如果各算法的迭代初值相同, 那么, 在一定假设下, 有关系式  $x_k^I = x_k^A$  和  $x_k^{II} = x_k^G$  成立, 即 AM 的中间步和对应的一类 Krylov 子空间方法的迭代步相同. Walker et al. (2011) 最早对此关系给出严格证明, 并在文中称之为“基本等价”, 本文沿用这样的说法. 简要说来, 在线性情形, 可以证明  $\text{range}(X_k)$  是 Krylov 子空间, AM 的投影条件实质对应于两类 Galerkin 条件 (Saad, 2003), 因此可以证明基本等价性. 这种等价性解释了 AM 在线性方程组求解中能快速收敛的原因, AM 也因此被称为非线性 Krylov 子空间法 (Calef et al., 2013). 同时, Walker et al. (2011) 也指出 AM 在线性方程组求解中不如 GMRES 可靠. 此外, 前述交替迭代方法也和 GMRES 在一定情形下有等价性 (Lupo Pasini, 2019).

AM 在非线性定点问题中的收敛性分析直到 2015 年才有实质的突破, 目前的主要工作集中在对第二类 AM 的分析上. Toth et al. (2015) 在一定的假设条件下证明了 AM-II( $m$ ) 的局部线性收敛性. 设问题 (1) 中  $g$  在定点的一个邻域内 Lipschitz 连续可微, Lipschitz 常数为  $\kappa \in (0, 1)$ . 对  $\tilde{\kappa} \in (\kappa, 1)$ , 如果初值足够好, 并且迭代中保证  $\sum_{j=0}^{m_k} |\alpha_k^{(j)}|$  一致有界, 那么对残差  $r_k$ , 有

$$\|r_k\|_2 \leq \tilde{\kappa}^k \|r_0\|_2. \quad (8)$$

之后, 类似的收敛性结论对于 EDIIS 也被证明成立 (Chen et al., 2019), 并于近期被推广到针对非光滑问题的分析中 (Mai et al., 2020; Bian et al., 2021; Bian et al., 2022). 然而, 这些理论结果还非常局限. Anderson (2019) 在其评论中指出, 结论 (8) 不能解释 AM 在实践中比定点迭代明显更好的收敛性, 因为后者  $q$ -线性收敛并且  $q$ -因子为  $\kappa$ . 为此, Evans et al. (2020) 给出了一个改进的结论:

$$\|r_{k+1}\|_2 \leq s_k (1 - \beta_k + \kappa \beta_k) \|r_k\|_2 + \sum_{j=0}^m \mathcal{O}(\|r_{k-j}\|_2^2), \quad (9)$$

其中  $s_k = \|\bar{r}_k\|_2 / \|r_k\|_2$ ,  $k \geq m$ . 因为总有  $s_k \in [0, 1]$ , 因此当  $s_k \ll 1$  时, 残差的一阶分量迅速衰减. Pollock et al. (2021) 给出了更精细的分析结果. 由于  $s_k$  是在迭代过程中才能确定的, 结论 (9) 不能在迭代之初预测 AM 的收敛情况.

## 2 正则化的 Anderson 混合

正则化的 Anderson 混合 (Scieur et al., 2020) 是 AM 的一种改进算法, 通常能够改善 AM 迭代的稳定性. 在 AM 的外推计算中, 求解最小二乘问题可能会带来数值稳定性问题, 因此 Walker et al. (2011) 建议使用 QR 分解求解问题 (7), 这样在  $R_k$  的列接近线性相关时可以确保求解的精度. 即便如此, 对于求解非线性问题, 如果算得的  $|\alpha_k^{(j)}|$  过大, 也会使得迭代不稳定. 因此, 正则化的 AM 在问题 (7) 中引入正则项以限制  $\|\Gamma_k\|_2$  的大小,  $\Gamma_k$  的确定方式为

$$\Gamma_k = \arg \min_{\Gamma \in \mathbb{R}^{m_k}} \|\mathbf{r}_k - \mathbf{R}_k \Gamma\|_2^2 + \delta \|\Gamma\|_2^2, \quad (10)$$

其中  $\delta \geq 0$  为正则化参数. 从而得到  $\Gamma_k = (\mathbf{R}_k^T \mathbf{R}_k + \delta \mathbf{I})^{-1} \mathbf{R}_k^T \mathbf{r}_k$ , 缓解了  $\mathbf{R}_k$  非列满秩时的数值影响.

正则化的 AM 是人们应用 AM 求解实际问题的一种常用方案. Scieur et al. (2020) 在交替迭代算法中引入该正则化, 得到适用于无约束最优化的优化算法, 并且通过正则化的 Chebyshev 多项式给出了收敛性分析. Fu et al. (2020) 将正则化的 AM 用于 Douglas-Rachford 算子分裂迭代的加速, 算法能够求解带线性等式约束的非光滑凸优化问题. Henderson et al. (2019) 在正则化的 AM 基础上引入重启和单调性检验, 将算法应用于机器学习中的 EM 算法的加速. Sun et al. (2021) 将正则化的 AM 用于强化学习的训练加速. 在这些实践中, 正则化对算法的稳定性起到了重要作用.

### 3 随机 Anderson 混合

随机 Anderson 混合 (Wei et al., 2021) 是 AM 的一种随机版本, 适用于求解随机优化问题. 问题描述为

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \mathbb{E}_\xi [F(\mathbf{x}; \xi)], \quad (11)$$

其中函数  $F: \mathbb{R}^d \times \Xi \rightarrow \mathbb{R}$  是连续可微并且可能非凸的函数,  $\xi \in \Xi$  是随机变量. 因为通常情况下  $F(\cdot; \xi)$  的具体形式难以显式给出, 如  $\xi$  服从一个未知的概率分布, 或者显式计算  $f$  开销过大, 所以实践中只能获得问题 (11) 的带噪声的一阶信息, 即带噪声的梯度估计. 通过对  $\xi$  采样, 得到问题 (11) 的一个特例, 即经验风险极小化问题:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{T} \sum_{i=1}^T f_i(\mathbf{x}), \quad (12)$$

其中  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  是关于第  $i$  个数据样本的函数,  $T$  是样本总数. 问题 (12) 广泛存在于机器学习的各种算法之中. 通常  $T$  很大, 造成遍历数据集得到全梯度  $\nabla f(\mathbf{x})$  的代价昂贵, 因此实用的方式是在  $T$  个样本中随机采样, 得到样本集上的梯度作为全梯度的估计.

使用 AM 求解优化问题是一个自然的想法, 因为梯度下降法  $\mathbf{x}_{k+1} = g(\mathbf{x}_k) := \mathbf{x}_k - \nabla f(\mathbf{x}_k)$  是一个定点迭代, 可以尝试用 AM 加速, 此时残差为  $\mathbf{r}_k = g(\mathbf{x}_k) - \mathbf{x}_k = -\nabla f(\mathbf{x}_k)$ . 然而, 随机优化有本质的难度, 由于不能得到精确的梯度, 如果使用带噪声的梯度定义残差, 传统的 AM 没有任何收敛性保证. 随机 AM 将 AM 推广到求解随机优化, 拓展了 AM 的应用范围.

在随机 AM 算法中, 定义残差  $\mathbf{r}_k = -g(\mathbf{x}_k)$ , 其中  $g(\mathbf{x}_k)$  是无偏的梯度估计. 相应地, 如式 (5) 所示可以得到历史序列  $\mathbf{X}_k$ ,  $\mathbf{R}_k \in \mathbb{R}^{d \times m_k}$ , 其中  $m_k = \min\{m, k\}$ . 随机 AM 在 AM-II( $m$ ) 的基础上引入了阻尼投影和自适应正则化, 其投影步和混合步为:

$$\bar{\mathbf{x}}_k = \mathbf{x}_k - \alpha_k \mathbf{X}_k \Gamma_k, \quad (\text{投影步}) \quad \bar{\mathbf{r}}_k = \mathbf{r}_k - \alpha_k \mathbf{R}_k \Gamma_k, \quad \mathbf{x}_{k+1} = \bar{\mathbf{x}}_k + \beta_k \bar{\mathbf{r}}_k, \quad (\text{混合步}) \quad (13)$$

其中  $\alpha_k \in [0, 1]$  为阻尼参数, 系数  $\Gamma_k$  由以下正则化的最小二乘问题确定:

$$\Gamma_k = \arg \min_{\Gamma \in \mathbb{R}^{m_k}} \|\mathbf{r}_k - \mathbf{R}_k \Gamma\|_2^2 + \delta_k \|\mathbf{X}_k \Gamma\|_2^2, \quad (14)$$

其中  $\delta_k \geq 0$  为正则化参数. 由于投影步的变化可能过大, 导致中间步越过目标函数当前的信赖域, 因此阻尼投影使得  $\bar{\mathbf{x}}_k = (1 - \alpha_k) \mathbf{x}_k + \alpha_k (\mathbf{x}_k - \mathbf{X}_k \Gamma_k) = \mathbf{x}_k - \alpha_k \mathbf{X}_k \Gamma_k$ . 同时, 正则化也起到限制  $\|\mathbf{X}_k \Gamma_k\|_2$  的大小的作用. 这些操作可以改善算法在随机场景中的鲁棒性和稳定性, 确保算法的收敛性.

随机 AM 在非凸随机优化中有全局收敛性. 假设  $f$  连续可微且有下界,  $\nabla f$  全局 Lipschitz 连续; 各样本独立无关, 并且与之前的迭代步无关, 梯度估计是全梯度的无偏估计, 且方差一致有界. 对于随机 AM, 使用递减的混合参数

$$\sum_{k=0}^{+\infty} \beta_k = +\infty, \quad \sum_{k=0}^{+\infty} \beta_k^2 < +\infty,$$

并对  $\alpha_k, \delta_k$  施加必要的条件, 有

$$\liminf_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\|_2 = 0 \quad \text{以概率 1 成立.}$$

如果还有梯度估计  $g(\mathbf{x}_k)$  一致有界, 那么有  $\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\|_2 = 0$  以概率 1 成立. 此外, 如果从历史迭代步

里随机选取解  $\bar{\mathbf{x}}$ , 那么为了确保  $\mathbb{E}\left[\|\nabla f(\bar{\mathbf{x}})\|_2^2\right] \leq \epsilon$ , 所需的梯度采样总次数为  $\mathcal{O}(\epsilon^{-2})$ , 这表明随机 AM 达到了一阶黑盒随机优化的最优复杂度.

此外, Wei et al.(2021)还给出随机 AM 的改进方案, 方案包括预处理和方差减少技术.

预处理的随机 AM 使用了预处理的混合步:

$$\mathbf{x}_{k+1} = \bar{\mathbf{x}}_k + \mathbf{M}_k^{-1} \bar{\mathbf{r}}_k, \quad (15)$$

其中  $\mathbf{M}_k$  是对 Hessian 阵  $\nabla^2 f(\mathbf{x}_k)$  的近似. 将式(13)中的投影步和式(15)结合即为预处理的随机 AM. 设整体的迭代式为  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{G}_k \mathbf{r}_k$ , 那么在  $\alpha_k = 1, \delta_k = 0$  时,  $\mathbf{G}_k$  求解了  $\min_{\mathbf{G}} \|\mathbf{G} - \mathbf{M}_k^{-1}\|_F$  s.t.  $\mathbf{G} \mathbf{R}_k = -\mathbf{X}_k$ .

方差减少技术起源于 SVRG 算法 (Johnson et al., 2013), 适用于求解经验风险极小化问题. 如果在随机 AM 中使用方差减少的梯度估计, 那么为了确保  $\mathbb{E}\left[\|\nabla f(\mathbf{x})\|_2^2\right] \leq \epsilon$ , 所需的梯度采样总次数为  $\mathcal{O}(T + (T^{2/3}/\epsilon))$ , 即算法复杂度得到了改进.

随机 AM 已经被成功用于深度学习的神经网络训练. 在图像分类和语言模型等任务上, 随机 AM 相比现有的随机优化器有明显更好的收敛性, 在大部分任务上节约了总的计算时间, 因此继承了 AM 在确定性问题中的优良效果, 在非凸随机优化中有很好的适用性.

## 4 短递归的 Anderson 混合

短递归的 Anderson 混合 (Wei et al., 2022a) 减少 AM 的内存开销, 适用于高维问题的求解. 与各种类型的拟 Newton 法类似, AM 需要存储历史序列  $\mathbf{X}_k, \mathbf{R}_k \in \mathbb{R}^{d \times m_k}$ , 因此相比定点迭代需要额外存储  $2m_k$  个维数为  $d$  的向量. 如果历史序列长度较大, 内存开销将成为 AM 的瓶颈, 致使算法无法在存储资源受限的机器上求解高维问题. 短递归的 AM 将历史序列的长度降到 2, 同时能够保证良好的收敛性.

首先介绍短递归 AM 的基础形式. 与 AM 不同, 短递归的 AM 使用修正的历史序列  $\mathbf{P}_k, \mathbf{Q}_k \in \mathbb{R}^{d \times 2}$ , 在每步迭代之初需要对向量对  $\Delta \mathbf{x}_{k-1}, \Delta \mathbf{r}_{k-1}$  作修正. 初始化  $\mathbf{P}_0, \mathbf{Q}_0 = \mathbf{0}$ , 在第  $k$  步迭代, 构造  $\mathbf{p}_k, \mathbf{q}_k \in \mathbb{R}^d$ :

$$\mathbf{p}_k = \Delta \mathbf{x}_{k-1} - \mathbf{P}_{k-1} \boldsymbol{\zeta}_k, \quad \mathbf{q}_k = \Delta \mathbf{r}_{k-1} - \mathbf{Q}_{k-1} \boldsymbol{\zeta}_k, \quad (16)$$

其中选取  $\boldsymbol{\zeta}_k = \arg \min_{\boldsymbol{\zeta}} \|\Delta \mathbf{r}_{k-1} - \mathbf{Q}_{k-1} \boldsymbol{\zeta}\|_2$ . 从而得到  $\mathbf{P}_k = (\mathbf{p}_{k-1}, \mathbf{p}_k), \mathbf{Q}_k = (\mathbf{q}_{k-1}, \mathbf{q}_k)$ . 进而迭代为

$$\bar{\mathbf{x}}_k = \mathbf{x}_k - \mathbf{P}_k \boldsymbol{\Gamma}_k \text{ (投影步)}, \quad \mathbf{x}_{k+1} = \bar{\mathbf{x}}_k + \beta_k \bar{\mathbf{r}}_k \text{ (混合步)}, \quad (17)$$

其中  $\bar{\mathbf{r}}_k = \mathbf{r}_k - \mathbf{Q}_k \boldsymbol{\Gamma}_k, \boldsymbol{\Gamma}_k = \arg \min_{\boldsymbol{\Gamma}} \|\mathbf{r}_k - \mathbf{Q}_k \boldsymbol{\Gamma}\|_2$ .

当求解对称正定线性方程组 (或强凸二次优化问题) 时, 由式(16)~(17)定义的短递归 AM 和全记忆的第二类 AM 完全等价, 这意味着短递归的 AM 虽然仅使用长度为 2 的历史序列, 但是却与 AM-II( $\infty$ ) 有相同的收敛性, 因此不存在历史信息的遗忘.

对于求解一般的非线性定点问题, 通过引入周期性重启和对  $\|\mathbf{P}_{k-1} \boldsymbol{\zeta}_k\|_2, \|\mathbf{Q}_{k-1} \boldsymbol{\zeta}_k\|_2$  的有界性检查, 在对  $g$  的标准假设 (Toth et al., 2015; Evans et al., 2020) 下, 短递归的 AM 具有局部线性收敛性:

$$\|\mathbf{r}_{k+1}\|_2 \leq s_k (|1 - \beta_k| + \kappa \beta_k) \|\mathbf{r}_k\|_2 + \hat{\kappa} \sum_{j=0}^{m_k} \mathcal{O}(\|\mathbf{r}_{k-j}\|_2^2), \quad (18)$$

其中  $s_k = \|\bar{\mathbf{r}}_k\|_2 / \|\mathbf{r}_k\|_2 \leq 1$ ,  $\kappa$  和  $\hat{\kappa}$  分别为  $g$  和  $g$  的导数的 Lipschitz 常数. 因此短递归 AM 在理论上没有减弱 AM 的收敛性. 对于求解非凸优化问题, 通过引入阻尼投影和正则化, 短递归的 AM 具有全局收敛性, 并且在非凸随机优化中有收敛性保证. 因此, 如果应用对迭代算法的内存占用有限制, 那么相较于有限内存的 AM 和随机 AM, 短递归的 AM 更有优势.

## 5 具有极小内存开销的 Anderson 混合

具有极小内存开销的 Anderson 混合 (Wei et al., 2022b) (Min-AM) 是一种基于 AM 的高效优化算法, 适用于大规模优化问题的求解. Min-AM 的历史序列长度为 1, 因此具有极小的内存开销. 同时, Min-AM 在优化问题中仍有不输于拟 Newton 法的收敛性.

对于求解优化问题, 因为光滑函数在最优解的局部邻域能用一个二次函数近似, 所以如果优化器能

快速地优化该二次函数, 那么在优化原目标函数时也有望有良好的收敛效果. 因此, 考虑强凸二次优化

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}, \quad (19)$$

其中  $\mathbf{A} \in \mathbb{R}^{d \times d}$  对称正定. 令  $\mathbf{s}_k := \Delta \mathbf{x}_{k-1}$ ,  $\mathbf{t}_k := \Delta \mathbf{r}_{k-1}$ , 那么 AM-II(1) 的迭代公式为  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{G}_k \mathbf{r}_k$ , 其中由式(6)~(7)得到  $\mathbf{G}_k = \beta_k \mathbf{I} - (\mathbf{s}_k + \beta_k \mathbf{t}_k)(\mathbf{t}_k^T \mathbf{t}_k)^{-1} \mathbf{t}_k^T$ . 从拟 Newton 法的观点来看,  $\mathbf{G}_k$  不是对称的, 因此不是一个对 Hessian 矩阵的逆矩阵的良好近似. Min-AM 使用修正的历史向量对  $\mathbf{p}_k, \mathbf{q}_k \in \mathbb{R}^d$ , 并且在 AM 的迭代格式基础上引入了附加的投影步. 具体描述如下:

初始化  $\mathbf{p}_1 = \mathbf{s}_0, \mathbf{q}_1 = \mathbf{t}_0$ . 对  $k \geq 2$ , 假设  $\mathbf{p}_{k-1}^T \mathbf{q}_{k-1} \neq 0$ , 那么在第  $k$  步迭代之初构造  $\mathbf{p}_k, \mathbf{q}_k$  为

$$\mathbf{p}_k = \mathbf{s}_k - \mathbf{p}_{k-1} \zeta_k, \quad \mathbf{q}_k = \mathbf{t}_k - \mathbf{q}_{k-1} \zeta_k,$$

其中  $\zeta_k = (\mathbf{p}_{k-1}^T \mathbf{q}_{k-1})^{-1} \mathbf{p}_{k-1}^T \mathbf{t}_k$ . 该操作使得  $\mathbf{q}_k \perp \mathbf{p}_{k-1}$ . 之后是对  $\mathbf{x}_k$  的迭代更新, 在式(6)之后引入附加的投影步, 即

$$\mathbf{x}_k^{(1)} = \mathbf{x}_k - \mathbf{p}_k \mathbf{\Gamma}_k^{(1)} \text{ (投影步)}, \quad \mathbf{x}_k^{(2)} = \mathbf{x}_k^{(1)} + \beta_k \mathbf{r}_k^{(1)} \text{ (混合步)}, \quad \mathbf{x}_{k+1} = \mathbf{x}_k^{(2)} - \mathbf{p}_k \mathbf{\Gamma}_k^{(2)} \text{ (投影步)}, \quad (20)$$

其中  $\mathbf{r}_k^{(1)} := \mathbf{r}_k - \mathbf{q}_k \mathbf{\Gamma}_k^{(1)}$ ,  $\beta_k > 0$ .  $\mathbf{x}_k^{(1)}, \mathbf{x}_k^{(2)}$  都是中间步, 为了确定  $\mathbf{\Gamma}_k^{(1)}, \mathbf{\Gamma}_k^{(2)}$ , 定义  $\mathbf{r}_k^{(2)} = \mathbf{r}_k^{(1)} - \beta_k \mathbf{A} \mathbf{r}_k^{(1)}$ , 可见  $\mathbf{r}_k^{(1)}, \mathbf{r}_k^{(2)}$  分别是  $\mathbf{x}_k^{(1)}, \mathbf{x}_k^{(2)}$  处的残差.  $\mathbf{\Gamma}_k^{(1)}, \mathbf{\Gamma}_k^{(2)}$  由第一类 AM 的投影条件确定:

$$\mathbf{r}_k^{(1)} = \mathbf{r}_k - \mathbf{q}_k \mathbf{\Gamma}_k^{(1)} \perp \mathbf{p}_k, \quad \mathbf{r}_{k+1} = \mathbf{r}_k^{(2)} - \mathbf{q}_k \mathbf{\Gamma}_k^{(2)} \perp \mathbf{p}_k. \quad (21)$$

假设  $\mathbf{p}_k^T \mathbf{q}_k \neq 0$ , 由式(20)给出的迭代为

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{H}_k \mathbf{r}_k, \quad \mathbf{H}_k = -\frac{\mathbf{p}_k \mathbf{p}_k^T}{\mathbf{p}_k^T \mathbf{q}_k} + \beta_k \left( \mathbf{I} - \frac{\mathbf{p}_k \mathbf{q}_k^T}{\mathbf{p}_k^T \mathbf{q}_k} \right) \left( \mathbf{I} - \frac{\mathbf{q}_k \mathbf{p}_k^T}{\mathbf{p}_k^T \mathbf{q}_k} \right). \quad (22)$$

可以看到  $\mathbf{H}_k$  是对称的. 迭代公式(22)即为 Min-AM 的基础形式.

在求解强凸二次优化问题(19)中, Wei et al.(2022b)揭示了 Min-AM 与共轭梯度法、Newton 法、BFGS(Nocedal et al., 2006)的本质联系. 以下介绍有关结论.

关于 Min-AM 的一个重要结论是 Min-AM 与第一类 AM 和共轭梯度法基本等价. 考虑求解问题(19), 设  $\{\mathbf{x}_k\}$  是 Min-AM 生成的迭代序列,  $\mathbf{x}_k^{(1)}$  是第  $k$  步迭代中的第一个中间步(见迭代格式(20)); 设  $\{\mathbf{x}_k^1\}$  是第一类 AM 生成的迭代序列,  $\bar{\mathbf{x}}_k^1$  是第  $k$  步迭代中的中间步(见迭代格式(6));  $\{\mathbf{x}_k^{\text{CG}}\}$  是共轭梯度法生成的迭代序列. 基本等价性指如果迭代初值相同, 那么

$$\mathbf{x}_k^{(1)} = \bar{\mathbf{x}}_k^1 = \mathbf{x}_k^{\text{CG}}. \quad (23)$$

这意味着 3 个算法的收敛性基本相同.

定义  $\mathbf{P}_k = (\mathbf{p}_1, \dots, \mathbf{p}_k)$ ,  $\mathbf{Q}_k = (\mathbf{q}_1, \dots, \mathbf{q}_k)$ , 并定义  $\mathbf{V}_k \in \mathbb{R}^{d \times (d-k)}$  使得  $\mathbf{V}_k^T \mathbf{P}_k = \mathbf{0}$ . 对优化问题(19), 在第  $k$  步迭代, 将  $\mathbf{x} \in \mathbb{R}^d$  写为  $\mathbf{x} = \mathbf{x}_k - \mathbf{P}_k \boldsymbol{\gamma} - \mathbf{V}_k \boldsymbol{\eta}$ , 其中  $\boldsymbol{\gamma} \in \mathbb{R}^k, \boldsymbol{\eta} \in \mathbb{R}^{d-k}$ . 先在子空间  $\text{range}(\mathbf{P}_k)$  上运用 Newton 法, 接着在  $\text{range}(\mathbf{P}_k)^\perp$  上以步长  $\beta_k$  梯度下降, 最后再在  $\text{range}(\mathbf{P}_k)$  上运用 Newton 法, 得到

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{H}_k^{\text{B}} \mathbf{r}_k, \quad \mathbf{H}_k^{\text{B}} = -\mathbf{P}_k (\mathbf{P}_k^T \mathbf{Q}_k)^{-1} \mathbf{P}_k^T + \beta_k (\mathbf{I} - \mathbf{P}_k (\mathbf{P}_k^T \mathbf{Q}_k)^{-1} \mathbf{Q}_k^T) (\mathbf{I} - \mathbf{Q}_k (\mathbf{P}_k^T \mathbf{Q}_k)^{-1} \mathbf{P}_k^T). \quad (24)$$

称迭代格式(24)为 B 迭代. 通过计算可以验证  $\mathbf{H}_k^{\text{B}}$  求解了  $\min_{\mathbf{H}} \|\mathbf{H} - \beta_k \mathbf{I}\|_{F(A)}$ , s.t.  $\mathbf{H} \mathbf{Q}_k = -\mathbf{P}_k, \mathbf{H} = \mathbf{H}^T$ . 进一步地, 由 Min-AM 的性质可以导出关系式(Wei et al., 2022b):

$$\mathbf{H}_k \mathbf{r}_k = \mathbf{H}_k^{\text{B}} \mathbf{r}_k. \quad (25)$$

这表明在求解强凸二次优化问题时, Min-AM 和 B 迭代完全等价. 而 B 迭代由 Newton 法和梯度下降法导出, 并且是一种对称化的多重割线拟 Newton 法, 当  $k = d$  时, B 迭代在全空间上使用 Newton 法. 这就建立了 Min-AM 和 Newton 法的联系, 可以认为 Min-AM 隐式地构造了 Hessian 阵的近似逆矩阵  $\mathbf{H}_k^{\text{B}}$ .

进一步地, 如果 Min-AM 和 B 迭代的参数  $\beta_k$  为常数  $\beta$ , 那么有  $\mathbf{H}_0^{\text{B}} = \beta \mathbf{I}$ , 随后的  $\mathbf{H}_k^{\text{B}}$  求解了

$$\min_{\mathbf{H}} \|\mathbf{H} - \mathbf{H}_{k-1}^{\text{B}}\|_{F(A)}, \quad \text{s.t. } \mathbf{H} \mathbf{q}_k = -\mathbf{p}_k, \mathbf{H} = \mathbf{H}^T. \quad (26)$$

如果将  $\mathbf{p}_k, \mathbf{q}_k$  替换为  $\mathbf{s}_k, \mathbf{t}_k$ , 那么式(26)导出 BFGS 算法. 这个关系表明, B 迭代使用修正的历史序列构造 Hessian 阵的近似逆矩阵, 由于 Min-AM 和 B 迭代等价, 因此 Min-AM 相较于 BFGS 能够减少大量的内存占用.

对于求解一般的非线性光滑优化乃至随机优化, Min-AM 也有明确的收敛性结论. 在确定性的光滑优

化中, 通过在迭代格式(22)上引入重启动和必要的检验, 可以证明 Min-AM 的收敛率最优地依赖于问题的条件数, Min-AM 和使用精确线搜索的非线性共轭梯度法有相当的收敛性. 在随机优化中, 与随机 AM 类似, 引入阻尼项和正则化, Min-AM 有全局收敛性并达到了最优的迭代复杂度. 因此, Min-AM 不仅将 AM 在优化问题中的内存开销降到极小, 而且保证了算法的收敛性, 在优化问题的求解中具备明显的优势. 此外, 对于确定性的光滑优化, Wei et al.(2022b)指出 Min-AM 可以使用很小的附加计算代价估计 Hessian 矩阵的特征值信息, 从而估计混合参数  $\beta_k$  的最优选取, 这对于算法的实际应用也是有益的.

## 6 总结

Anderson 混合是加速定点迭代的一种强有力的算法, 现有的研究揭示了其与 Krylov 子空间方法和拟 Newton 法的深刻联系. 目前 Anderson 混合的收敛性问题还没有得到完全解决, 仍需要有更好的理论分析结果来更精确地刻画 Anderson 混合在非线性定点问题中的收敛行为. 为了改善算法的稳定性、增大算法的使用范围, 一些 Anderson 混合的改进算法被提出并得到了成功应用. 本文介绍了其中一些有代表性的改进算法, 结论表明基于 Anderson 混合的新算法能够被用于随机优化等更困难的问题, 并且能够在内存开销上相较于传统的拟 Newton 法有明显的优势, 有望解决科学计算和机器学习等领域中具有挑战性的实际问题, 值得被进一步研究.

### 参考文献:

- ANDERSON D G, 1965. Iterative procedures for nonlinear integral equations[J]. J ACM, 12(4): 547-560.
- ANDERSON D G, 2019. Comments on "Anderson acceleration, mixing and extrapolation" [J]. Numer Algorithms, 80(1): 135-234.
- ARORA A, MORSE D C, BATES F S, et al, 2017. Accelerating self-consistent field theory of block polymers in a variable unit cell[J]. J Chem Phys, 146(24): 244902.
- BIAN W, CHEN X J, 2022. Anderson acceleration for nonsmooth fixed point problems[J]. SIAM J Numer Anal, 60(5): 2565-2591.
- BIAN W, CHEN X J, KELLEY C T, 2021. Anderson acceleration for a class of nonsmooth fixed-point problems[J]. SIAM J Sci Comput, 43(5): S1-S20.
- BREZINSKI C, REDIVO-ZAGLIA M, SAAD Y, 2018. Shanks sequence transformations and Anderson acceleration [J]. SIAM Rev, 60(3): 646-669.
- CALEF M T, FICHTL E D, WARSA J S, et al, 2013. Nonlinear Krylov acceleration applied to a discrete ordinates formulation of the  $k$ -eigenvalue problem[J]. J Comput Phys, 238: 188-209.
- CHEN X J, KELLEY C T, 2019. Convergence of the EDIIS algorithm for nonlinear equations[J]. SIAM J Sci Comput, 41(1): A365-A379.
- EVANS C, POLLOCK S, REBHOLZ L G, et al, 2020. A proof that Anderson acceleration improves the convergence rate in linearly converging fixed-point methods (but not in those converging quadratically)[J]. SIAM J Numer Anal, 58(1): 788-810.
- FANG H R, SAAD Y, 2009. Two classes of multisection methods for nonlinear acceleration [J]. Numer Linear Algebra Appl, 16(3): 197-221.
- FU A Q, ZHANG J Z, BOYD S, 2020. Anderson accelerated Douglas-Rachford splitting[J]. SIAM J Sci Comput, 42(6): A3560-A3583.
- HENDERSON N C, VARADHAN R, 2019. Damped Anderson acceleration with restarts and monotonicity control for accelerating EM and EM-like algorithms[J]. J Comput Graph Stat, 28(4): 834-846.
- JOHNSON R, ZHANG T, 2013. Accelerating stochastic gradient descent using predictive variance reduction [C]//Proceedings of the 27th International Conference on Neural Information Processing Systems: 315-323.
- KUDIN K N, SCUSERIA G E, CANCÈS E, 2002. A black-box self-consistent field convergence algorithm: One step closer[J]. J Chem Phys, 116(19): 8255-8261.
- LUPO PASINI M, 2019. Convergence analysis of Anderson-type acceleration of Richardson's iteration[J]. Numer Linear Algebra

- Appl, 26(4): e2241.
- MAI V V, JOHANSSON M, 2020. Anderson acceleration of proximal gradient methods[C]//Proceedings of the 37th International Conference on Machine Learning: 6620–6629.
- NOCEDAL J, WRIGHT S J, 2006. Numerical optimization[M]. 2nd ed. New York: Springer.
- POLLOCK S, REBHOLZ L G, 2021. Anderson acceleration for contractive and noncontractive operators[J]. IMA J Numer Anal, 41(4): 2841–2872.
- POLLOCK S, REBHOLZ L G, XIAO M Y, 2019. Anderson-accelerated convergence of Picard iterations for incompressible Navier-Stokes equations[J]. SIAM J Numer Anal, 57(2): 615–637.
- PRATAPA P P, SURYANARAYANA P, PASK J E, 2016. Anderson acceleration of the Jacobi iterative method: An efficient alternative to Krylov methods for large, sparse linear systems[J]. J Comput Phys, 306: 43–54.
- PULAY P, 1980. Convergence acceleration of iterative sequences. the case of SCF iteration[J]. Chem Phys Lett, 73(2): 393–398.
- ROHWEDDER T, SCHNEIDER R, 2011. An analysis for the DIIS acceleration method used in quantum chemistry calculations [J]. J Math Chem, 49(9): 1889–1914.
- SAAD Y, 1981. Krylov subspace methods for solving large unsymmetric linear systems[J]. Math Comp, 37(155): 105–126.
- SAAD Y, 2003. Iterative methods for sparse linear systems[M]. 2nd ed. Philadelphia: SIAM.
- SAAD Y, SCHULTZ M H, 1986. GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems[J]. SIAM J Sci Stat Comput, 7(3): 856–869.
- SCIEUR D, D'ASPROMONT A, BACH F, 2020. Regularized nonlinear acceleration[J]. Math Program, 179(1/2): 47–83.
- SUN K, WANG Y F, LIU Y, et al, 2021. Damped Anderson mixing for deep reinforcement learning: Acceleration, convergence, and stabilization[C]//Proceedings of the 35th Conference on Neural Information Processing Systems: 3732–3743.
- SURYANARAYANA P, PRATAPA P P, PASK J E, 2019. Alternating Anderson-Richardson method: An efficient alternative to preconditioned Krylov methods for large, sparse linear systems[J]. Comput Phys Commun, 234: 278–285.
- TOTH A, KELLEY C T, 2015. Convergence analysis for Anderson acceleration[J]. SIAM J Numer Anal, 53(2): 805–819.
- WALKER H F, NI P, 2011. Anderson acceleration for fixed-point iterations[J]. SIAM J Numer Anal, 49(4): 1715–1735.
- WEI F C, BAO C L, LIU Y, 2021. Stochastic Anderson mixing for nonconvex stochastic optimization [C]//Proceedings of the 35th Conference on Neural Information Processing Systems: 22995–23008.
- WEI F C, BAO C L, LIU Y, 2022a. A class of short-term recurrence Anderson mixing methods and their applications[C]//The Tenth International Conference on Learning Representations: 1565.
- WEI F C, BAO C L, LIU Y, et al, 2022b. A variant of Anderson mixing with minimal memory size[C]//Proceedings of the 36th Conference on Neural Information Processing Systems: 16650–16663.
- YANG Y N, 2021. Anderson acceleration for seismic inversion[J]. Geophysics, 86(1): R99–R108.

(责任编辑 冯兆永)