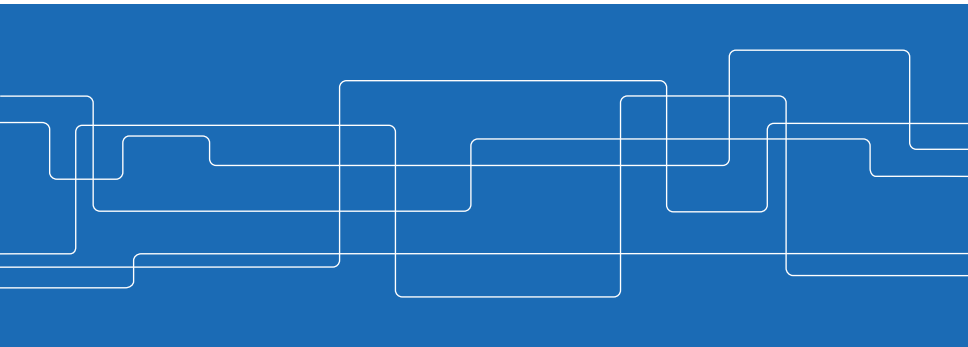




Anderson Acceleration of Constrained Optimization Algorithms

Vien V. Mai and Mikael Johansson
KTH - Royal Institute of Technology



Introduction

Generic **unconstrained** convex optimization problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} f(x)$$

Optimization algorithm produces **sequence** of iterates x_k converging to x^*

Example: gradient descent:

$$x_{k+1} = x_k - \gamma \nabla f(x_k)$$

Example: Newton's method:

$$x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$$

Note. Algorithms only keep the last iterate for the next update

Multi-step methods

Use multiple past iterates to accelerate the convergence process

Example: Heavy-ball method [Polyak, 1964]:

$$x_{k+1} = x_k + \beta_k(x_k - x_{k-1}) - \gamma \nabla f(x_k)$$

Example: Accelerated gradient descent [Nesterov, 1983]:

$$x_{k+1} = x_k + \beta_k(x_k - x_{k-1}) - \gamma \nabla f(x_k + \beta_k(x_k - x_{k-1}))$$

Optimal **momentum** parameters β_k often depend on unknown constants

Anderson acceleration [Anderson, 1965]

Originally developed for solving nonlinear integral equations

Recently adapted for solving general fixed-point equations

$$g(x) = x$$

Key ideas. Make clever use of past iterates

- Keeps $m + 1$ most recent iterates
- **Ideally**, forms $x_{\text{ext}} = \sum_{i=0}^m \alpha_i x_{k-i}$ such that

$$\alpha^* = \underset{\alpha: \alpha^\top \mathbf{1} = 1}{\operatorname{argmin}} \|g(x_{\text{ext}}) - x_{\text{ext}}\|.$$

- Finds coefficients α_i that minimize

$$\alpha \leftarrow \underset{\alpha^\top \mathbf{1} = 1}{\operatorname{argmin}} \left\| \sum_{i=0}^m \alpha_i (g(x_{k-i}) - x_{k-i}) \right\|$$

- Sets $x_{k+1} = \sum_{i=0}^m \alpha_i x_{k-i}$

Optimization algorithms as fixed-point iterations

Many optimization methods can be written as fixed-point iterations

Gradient descent:

$$x_{k+1} = x_k - \gamma \nabla f(x_k) \quad \Leftrightarrow \quad x_{k+1} = g(x_k)$$

with $g(x) := x - \gamma \nabla f(x)$.

Finding an optimal solution reduces to finding a fixed-point of g :

$$g(x^*) = x^* \quad \Leftrightarrow \quad \nabla f(x^*) = 0 \quad \Leftrightarrow \quad x^* \in \operatorname{argmin} f(x)$$

Idea. Apply Anderson acceleration to speed-up optimization algorithms

[Scieur, d'Aspremont, Bach, 2016], [Zhang, O'Donoghue, and Boyd, 2018]

Anderson acceleration for gradient descent

Goal. Find a point x^* satisfying $\nabla f(x^*) = 0$

AA-GD.

- Keeps $m + 1$ most recent iterates
- Finds coefficients α_i such that

$$\alpha \leftarrow \operatorname{argmin}_{\alpha^\top \mathbf{1} = 1} \left\| \sum_{i=0}^m \alpha_i \nabla f(x_{k-i}) \right\|$$

- Sets $x_{k+1} \leftarrow \sum_{i=0}^m \alpha_i x_{k-i}$

If f is convex quadratic and $m = k$, AA-GD is GMRES

Convergence rates for the quadratic objective

Accelerated gradient descent:

$$f(x_k) - f(x^\star) \leq O(1/k^2) (f(x_0) - f(x^\star))$$

Strong **practical performance** requires local adaption

- efficient line-search procedures [Nesterov, 2007]
- adaptive restart techniques [O'Donoghue and Candes, 2014]

Anderson acceleration for GD:

$$f(x_k) - f(x^\star) \leq O\left(\min\left\{1/k^2, e^{-\frac{k}{\sqrt{\kappa}}}\right\}\right) (f(x_0) - f(x^\star))$$

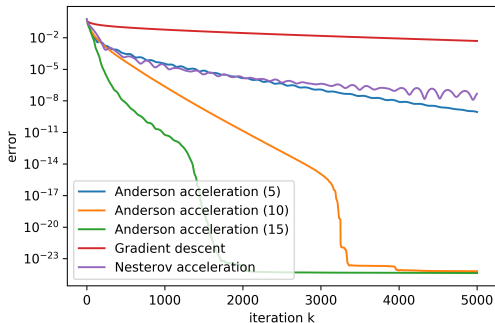
Very strong adaptive rate, remarkable property of Krylov subspace methods

Numerical example

Quadratic convex minimization

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} x^\top A x + b^\top x$$

- $A \in \mathbb{R}^{200 \times 200}$ with $\lambda_{\max}/\lambda_{\min} = 10^4$



Open problem: AA for constrained optimization

Consider **constrained** convex optimization problems

$$\underset{x \in \mathcal{C}}{\text{minimize}} f(x) \quad (1)$$

where \mathcal{C} is a closed convex set.

Recall the Anderson acceleration update

$$x_{k+1} = \sum_{i=0}^m \alpha_i x_{k-i}, \quad \alpha_i \in \mathbb{R}$$

The extrapolated point x_{k+1} may become **infeasible**

Q. Can we use Anderson acceleration for constrained problems?

Outline

- Background and motivation
- Accelerating projected gradient descent
- Local convergence of AA for constrained problems
- Numerical examples
- Summary and conclusions

Projected gradient descent

x^* is an optimal solution to (1) if and only if

$$x^* = \Pi_{\mathcal{C}}(x^* - \gamma \nabla f(x^*))$$

Projected gradient descent (PGD)

$$x_{k+1} = \Pi_{\mathcal{C}}(x_k - \gamma \nabla f(x_k))$$

Exactly the fixed-point iteration of the mapping

$$g(x) = \Pi_{\mathcal{C}}(x - \gamma \nabla f(x))$$

Naively using AA for g leads to iterates **infeasible**

Projected gradient descent

x^* is an optimal solution to (1) if and only if

$$x^* = \Pi_C(x^* - \gamma \nabla f(x^*))$$

An alternative representation of PGD:

$$y_{k+1} = x_k - \gamma \nabla f(x_k)$$

$$x_{k+1} = \Pi_C(y_{k+1})$$

Can be seen as the fixed-point iteration of the mapping

$$g(y) = \Pi_C(y) - \gamma \nabla f(\Pi_C(y))$$

Observation. x^* is optimal if and only if

$$y^* = g(y^*) \quad \text{and} \quad x^* = \Pi_C(y^*)$$

Anderson acceleration for PGD

Key idea. Use AA to speed-up fixed-point computations of $g(y)$

No iterates infeasibility since the sequence $\{y_k\}$ has no restriction

How to relate the convergence of $\{x_k\}$ and $\{y_k\}$:

$$\begin{aligned}\|x_{k+1} - x^*\| &= \|\Pi_C(y_{k+1}) - \Pi_C(x^* - \gamma \nabla f(x^*))\| \\ &\leq \|y_{k+1} - y^*\|\end{aligned}$$

► if AA quickly drives y_k to y^* , so does x_k to x^*

Convergence guarantee for general smooth mappings

Theorem 1 [Toth and Kelley, 2015]. Suppose that g is **differentiable** and contractive with constant ρ . If x_0 is sufficiently close to x^* , then

$$\|x_k - x^*\| \leq \rho^k \|x_0 - x^*\|.$$

So far, all convergence guarantees for AA rely on linearizing g around x^* :

$$g(x) = g(x^*) + G'(x^*) (x - x^*) + o(\|x - x^*\|)$$

Due to $\Pi_C(\cdot)$, the mapping $g(y)$ defined above is **non-differentiable**

Extending the analysis to non-smooth mappings

Let $F(x) \triangleq x - g(x)$, Theorem 1 indeed only needs the bounds

$$\|F(x) - F'(x^*)(x - x^*)\| \leq \frac{c}{2} \|x - x^*\|^2$$

for some constant $c > 0$ and for all x sufficiently close to x^* .

Extending Theorem 1 to **non-smooth** case amounts to searching for such F'

Two key ingredients:

- Clarke's generalized Jacobian
- (strong) **semi-smoothness**

Performance guarantees of AA-PGD

Lemma 1. Projections onto the nonnegative orthant, second-order cone, positive semidefinite cone, and polyhedral set are all strongly semi-smooth.

Main result. Let f be a μ -strongly convex and L -smooth function. Suppose that $\Pi_C(\cdot)$ is strongly semi-smooth and that x_0 is sufficiently close to x^* . Then,

$$\|x_k - x^*\| \leq \rho^k \|x_0 - x^*\|,$$

where $\rho = \sqrt{1 - \gamma 2\mu L / (\mu + L)}$.

Numerical results

Constrained logistic regression.

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{M} \sum_{i=1}^M \log(1 + \exp(-y_i a_i^\top x)) + \mu \|x\|^2 \\ & \text{subject to} \quad \|x\|_\infty \leq 1, \end{aligned}$$

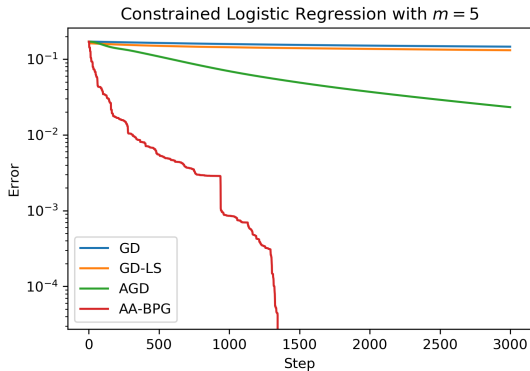
$a_i \in \mathbb{R}^n$ are training samples and $y \in \{-1, 1\}^n$ are the corresponding labels.

We use the UCI Madelon dataset with $M = 2000$ and $n = 500$

We set $\mu = 0.01$ and $m = 5$

➤ Extremely ill-conditioned problem with condition number $\kappa = 3 \times 10^9$

Numerical results



Numerical results

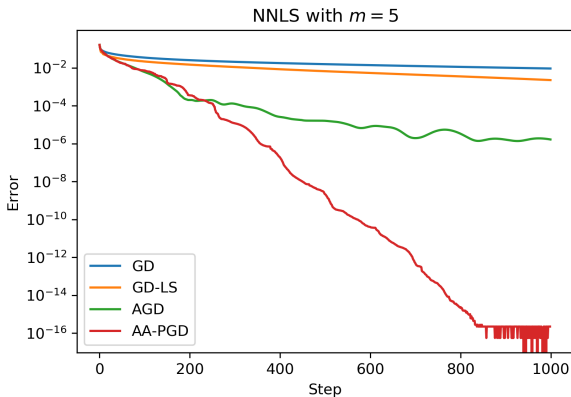


Figure: Nonnegative least-squares, $A \in \mathbb{R}^{200 \times 200}$ with $\lambda_{\max}/\lambda_{\min} = 10^4$

Conclusion

Anderson acceleration

- dramatic speed-ups in local convergence, at small extra cost
- current theory only applies to unconstrained problems

Our contributions

- first convergence results for AA on constrained problems
- strong practical performance

Future work

- algorithms: primal-dual methods, ADMM
- applications: Sinkhorn-Knopp, generative adversarial network (GAN)