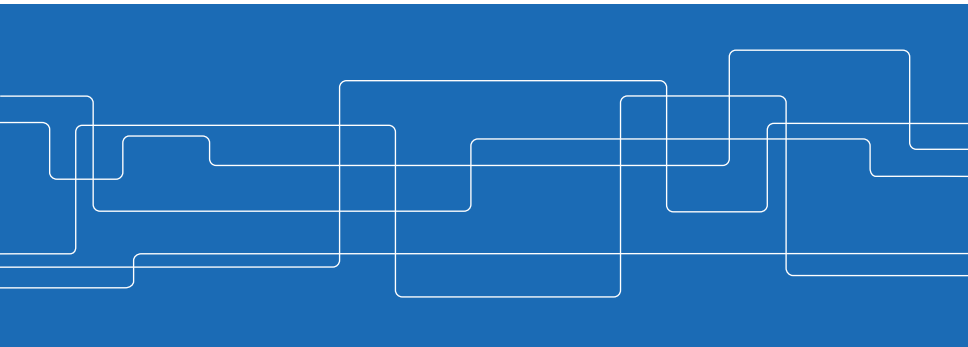




# Anderson Acceleration of Proximal Gradient Methods

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  $x^\star$

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

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

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

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

## 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.

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

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

## Convergence rates for the quadratic objective

Anderson acceleration for GD (AA-GD):

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

Very strong **adaptive** rate, remarkable property of Krylov subspace methods

Accelerated gradient descent:

$$f(x_k) - f(x^\star) \leq O\left(1/\textcolor{red}{k}^2\right) (f(x_0) - f(x^\star))$$

Strong **practical performance** requires local adaption

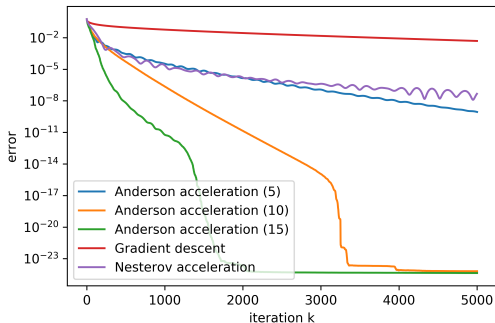
- adaptive restart techniques [O'Donoghue and Candes, 2014]

## 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

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

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

**Q.** Can we use Anderson acceleration for constrained problems?

## Outline

---

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

## Composite minimization and proximal operator

Consider composite convex optimization problems

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

- $f$  is smooth and convex;  $h$  is closed and convex

Proximal operator:

$$\text{prox}_h(y) := \underset{x \in \mathbb{R}^n}{\text{argmin}} \left\{ \frac{1}{2} \|x - y\|_2^2 + h(x) \right\}.$$

$x^*$  is an optimal solution if and only if

$$x^* = \text{prox}_{\gamma h}(x^* - \gamma \nabla f(x^*))$$

## Proximal gradient descent

Proximal gradient algorithm (PGA)

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

Exactly the fixed-point iteration of the mapping

$$g(x) = \text{prox}_{\gamma h}(x - \gamma \nabla f(x))$$

Finding  $x^*$  is equivalently to finding a fixed-point of  $g$

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

## Proximal gradient descent

An alternative representation of PGA:

$$\begin{aligned}y_{k+1} &= x_k - \gamma \nabla f(x_k) \\x_{k+1} &= \text{prox}_{\gamma h}(y_{k+1})\end{aligned}$$

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

$$g(y) = \text{prox}_{\gamma h}(y) - \gamma \nabla f(\text{prox}_{\gamma h}(y))$$

**Observation.**  $x^\star$  is optimal if and only if

$$y^\star = g(y^\star) \quad \text{and} \quad x^\star = \text{prox}_{\gamma h}(y^\star)$$



## 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  $\{y_k\}$ :

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

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

## Convergence guarantee

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

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

Note. The mapping  $g$  is non-differentiable

**Assumption.** The operator  $\text{prox}_h(\cdot)$  is **strongly semi-smooth**<sup>1</sup>

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

$$\|x_k - x^*\|^2 \leq \left(1 - \frac{\mu}{L}\right)^k \|x_0 - x^*\|^2.$$

---

<sup>1</sup>[Mifflin, 1977], [Qi and Sun, 1993]

## Numerical results

### Constrained logistic regression.

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

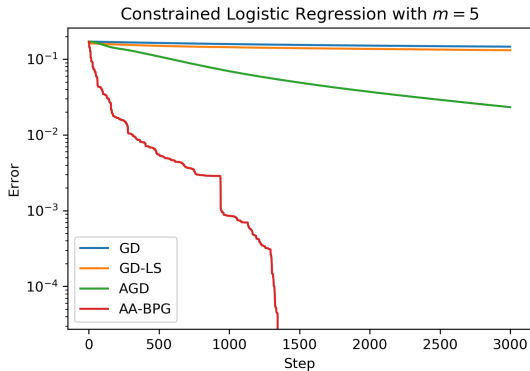
$a_i \in \mathbb{R}^n$  are training samples and  $y \in \{-1, 1\}^n$  are the corresponding labels.

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

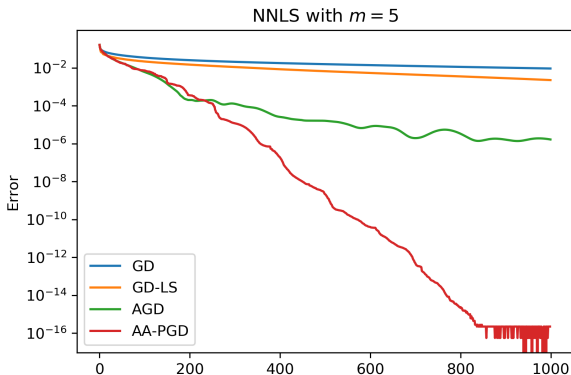


## Numerical results

### Nonnegative least squares.

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \|Ax - b\|^2 \quad \text{subject to } x \geq 0,$$

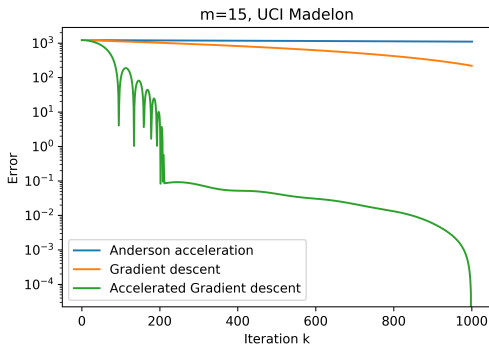
with  $A \in \mathbb{R}^{1000 \times 5000}$  and  $b \in \mathbb{R}^{1000}$ .



## Challenges with global convergence

Even for smooth problems, AA may not converge!

**Example.** AA-GD gets stuck in unconstrained logistic regression



## The real algorithm (Guarded-AA)

1. Computes

$$g_k = x_k - \gamma \nabla f(x_k) \quad (\text{gradient step})$$

2. Applies AA-PGA

$$y_{\text{ext}} = \sum_{i=0}^{m_k} \alpha_i^k g_{k-i} \quad \text{and} \quad x_{\text{test}} = \text{prox}_{\gamma h}(y_{\text{ext}})$$

3. Guarded step

$$\text{If } f(x_{\text{test}}) \leq f(x_k) - \frac{\gamma}{2} \|\nabla f(x_k)\|_2^2$$

$$x_{k+1} = x_{\text{test}}, \quad y_{k+1} = y_{\text{ext}}$$

**else**

$$x_{k+1} = \text{prox}_{\gamma h}(g_k), \quad y_{k+1} = g_k$$

**end**

## Numerical results

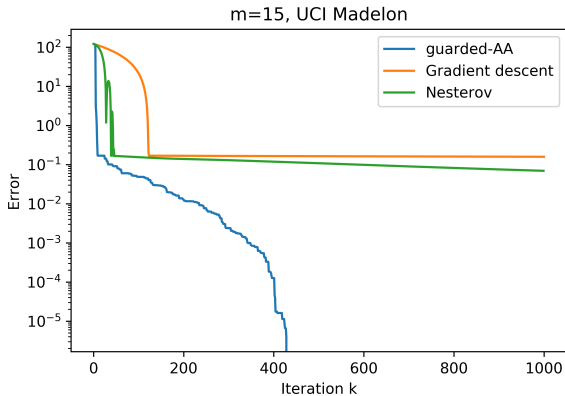


Figure: Bounded logistic regression



## Extending to Non-Euclidean Geometry

Consider optimization problems of the form

$$\underset{x \in \mathcal{D}}{\text{minimize}} \quad f(x) + h(x).$$

Problem geometry is exploited by a **kernel** function  $\varphi$

**Example.** Energy function  $\varphi(x) = (1/2) \|x\|_2^2$

**Example.** Shannon entropy  $\varphi(x) = \sum_{i=1}^n x_i \log x_i$ ,  $\text{dom } \varphi = \mathbb{R}_+^n$

Bregman proximal operator:

$$\text{prox}_h^\varphi(y) = \underset{x \in \mathbb{R}^n}{\text{argmin}} \{h(x) + D_\varphi(x, y)\}, \quad y \in \text{int dom } \varphi.$$

## Bregman proximal gradient (BPG)

Bregman proximal gradient (BPG)

$$x_{k+1} = \operatorname{argmin}_{x \in \mathbb{R}^n} \{ \langle \nabla f(x_k), x - x_k \rangle + \gamma^{-1} D_\varphi(x, x_k) + h(x) \}.$$

Can be expressed as

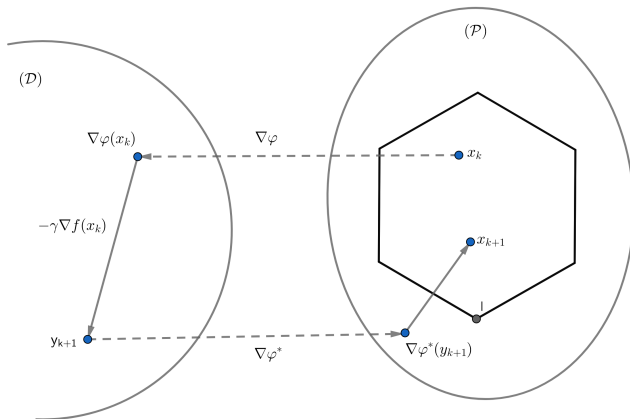
$$x_{k+1} = \operatorname{prox}_{\gamma h}^\varphi (\nabla \varphi^* (\nabla \varphi(x_k) - \gamma \nabla f(x_k)))$$

Equivalent form

$$\begin{aligned} y_{k+1} &= \nabla \varphi(x_k) - \gamma \nabla f(x_k) \\ x_{k+1} &= \operatorname{prox}_{\gamma h}^\varphi (\nabla \varphi^*(y_{k+1})). \end{aligned}$$

Mirror Descent is a special instance

# Illustration of Bregman Proximal gradient



## Anderson acceleration of BPG

**Strategy:** Extrapolate the **dual** sequence  $\{y_k\}$ :

$$g(y) = \nabla\varphi(\text{prox}_{\gamma h}^{\varphi} \circ \nabla\varphi^*(y)) - \gamma\nabla f(\text{prox}_{\gamma h}^{\varphi} \circ \nabla\varphi^*(y)).$$

**Assumption.** The conjugate of  $\varphi$  has full domain, i.e.,  $\text{dom } \varphi^* = \mathbb{R}^n$ .

Reduced to AA-PGA if  $\varphi(\cdot) = (1/2) \|\cdot\|_2^2$

A similar guarded step as in AA-PGA guarantees global convergence

## Relative-entropy nonnegative regression

The task is to reconstruct the signal  $x \in \mathbb{R}_+^n$  by solving

$$\underset{x}{\text{minimize}} D_{\text{KL}}(Ax, b) + \lambda \|x\|_1 \quad \text{subject to } x \geq 0,$$

- $A \in \mathbb{R}_+^{m \times n}$  is given nonnegative observation matrix
- $b \in \mathbb{R}_+^m$  is a noisy measurement vector

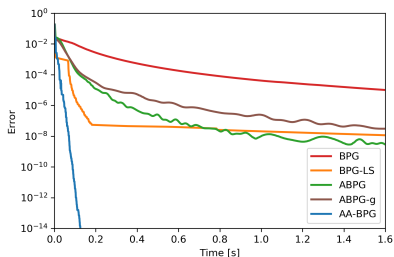
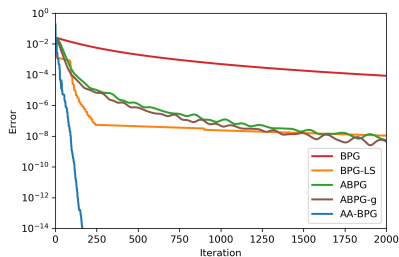
We adapt the family of BPG methods with:

- $\mathcal{D} = \mathbb{R}_+^n$
- $\varphi$  is the Shannon entropy,  $f(x) = D_{\text{KL}}(Ax, b)$
- $h(x) = \lambda \|x\|_1$  with  $\lambda = 0.001$

Compare with BPG, BPG-LS, ABPG, ABPG-LS

[Bauschke, Bolte, Teboulle, 2016], [Hanzely, Richtárik, Xiao, 2018]

## Numerical results



(a)  $(m,n) = (1000,100)$

## 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
- first application to non-Euclidean geometry

### Future work

- algorithms: primal-dual methods, ADMM
- applications: Sinkhorn-Knopp, optimal transport