

Accelerated Proximal Algorithms

Karol Chojnacki

March 2025

Gradient Descent (GD)

Optimization Problem:

$$\min_x f(x)$$

Update Rule:

$$x_{t+1} = x_t - \eta \nabla f(x_t)$$

Convergence Rate: $\mathcal{O}(1/t)$

Accelerated Gradient Descent (AGD)

Optimization Problem:

$$\min_x f(x)$$

Update Rule:

$$y_{t+1} = x_t - \eta \nabla f(x_t)$$

$$x_{t+1} = y_{t+1} + \frac{k-1}{k+2}(y_{t+1} - y_t)$$

Convergence Rate: $\mathcal{O}(1/t^2)$

Proximal Operator: For a closed, proper, convex function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$, the (scaled) proximal operator is defined as:

$$\text{prox}_{\lambda f}(v) = \arg \min_x \left(f(x) + \frac{1}{2\lambda} \|x - v\|^2 \right),$$

where $\lambda > 0$ scales the problem.

Motivation.

When we linearize $f(x)$ by $\langle \mathbf{x} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle$, prox iterating becomes GD.

$$\mathbf{x}^* = \operatorname{argmin} f(x) \iff \mathbf{x}^* = \operatorname{argmin} f(x) + \delta \|\mathbf{x} - \mathbf{x}^*\|_2^2$$

$$\begin{aligned} \mathbf{x}_{k+1} &= \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^N} \left(f(\mathbf{x}_k) + \langle \mathbf{x} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right) \\ &= \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^N} \left(\frac{\alpha_k}{2} \|\nabla f(\mathbf{x}_k)\|_2^2 + \langle \mathbf{x} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right) \\ &= \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^N} \left(\frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_k + \alpha_k \nabla f(\mathbf{x}_k)\|_2^2 \right) \\ &= \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \end{aligned}$$

Optimization Problem:

$$\min_{\mathbf{x}} g(\mathbf{x}) + h(\mathbf{x})$$

where $h(\mathbf{x})$ is a regularization term, e.g., $h(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$.

$$\begin{aligned}\mathbf{x}_{k+1} &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(g(\mathbf{x}_k) + \langle \mathbf{x} - \mathbf{x}_k, \nabla g(\mathbf{x}_k) \rangle + h(\mathbf{x}) + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right) \\ &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(h(\mathbf{x}) + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_k + \alpha_k \nabla g(\mathbf{x}_k)\|_2^2 \right) \\ &= \text{prox}_{\alpha_k h}(\mathbf{x}_k - \alpha_k \nabla g(\mathbf{x}_k)).\end{aligned}$$

Update Rule:

$$x_{t+1} = \text{prox}_{\lambda\eta h}(x_t - \eta\nabla g(x_t)) = T(x_t - \eta\nabla g(x_t))$$

Convergence Rate: $\mathcal{O}(1/t)$

Proximal Operator for L1 Regularization:

$$\text{prox}_{\lambda\|x\|_1}(v) := T(v) = \begin{cases} v - \lambda, & v > \lambda \\ 0, & |v| \leq \lambda \\ v + \lambda, & v < -\lambda \end{cases}$$

Optimization Problem:

$$\min_x f(x) + g(x)$$

where $g(x)$ is a regularization term.

Update Rule:

$$y_{t+1} = \text{prox}_{\lambda\eta g}(x_t - \eta\nabla f(x_t))$$

$$x_{t+1} = y_{t+1} + \frac{t-1}{t+2}(y_{t+1} - y_t)$$

Convergence Rate: $\mathcal{O}(1/t^2)$

Modification: Adaptive step size $\eta \rightarrow \eta_k$:

$$\eta_k = \frac{1}{L} + k\mu$$

Expected Convergence: Empirically close to $\mathcal{O}(1/t^2)$

Modification: Step-size depends on support structure:

$$\eta_k = c \frac{\|s_k \nabla f(x_k)\|^2}{\|\Phi(s_k \nabla f(x_k))\|^2}$$

where:

- c is a scaling constant controlling the step size.
- s_k is a binary mask indicating which elements of x_k are nonzero (support of x_k).
- $\nabla f(x_k)$ is the gradient of the loss function with respect to the model parameters.
- $\Phi(\cdot)$ represents a transformation that captures structural information, such as a convolutional layer in a CNN.

Expected Convergence: Adaptive, may accelerate convergence in sparse problems to $\mathcal{O}(1/t^3)$, *Gustavo Silva and Paul Rodriguez*.

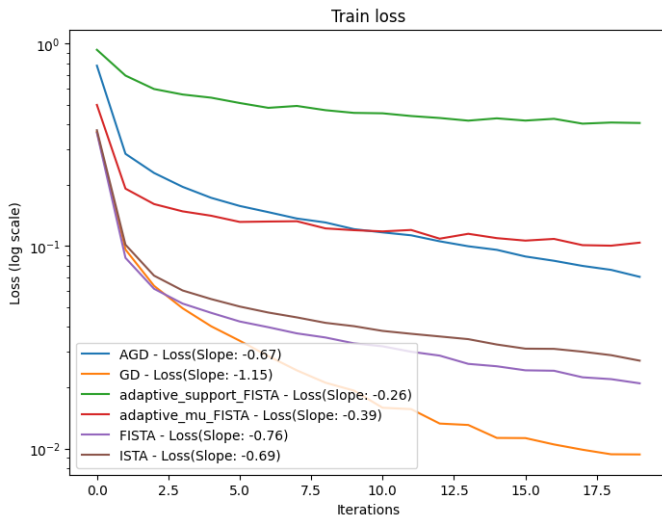
Goal: Compare optimization methods on ML problems.

Tasks: MNIST classification with CNN.

Loss function: Cross-entropy + L1.

Metrics: Convergence speed, sparsity, test accuracy.

Implementation: PyTorch, Google Colab (T4 GPU).



Algorithm	Final Sparsity (%)	Test Accuracy (%)
AGD	0.064	97.73
GD	0.075	98.25
adaptive_support_FISTA	50.968	63.10
adaptive_mu_FISTA	13.115	94.81
FISTA	39.842	98.31
ISTA	43.394	98.24

- Proximal Algorithms - Stanford University, Neal Parikh
- <https://sci-hub.se/https://ieeexplore.ieee.org/document/8903154>
- My code for experiment, if anybody wants hit me up

Thank you very much for listening!