

Proximal Gradient Algorithms

Proximal algorithms are particularly useful when the functional we are minimizing can be broken into two parts, one of which is smooth, and the other for which there is a fast proximal operator.

We are interested in minimizing unconstrained programs of the form

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x}),$$

where

1. $g(\mathbf{x})$ is convex and differentiable, and
2. $h(\mathbf{x})$ is convex with an inexpensive prox operator (examples to come).

The **proximal gradient** algorithm minimizes f iteratively, with each iteration consisting of

1. taking a gradient step on g , then
2. taking a proximal step on h .

We can combine these steps into one expression:

$$\mathbf{x}^{(k+1)} = \text{prox}_{t_k h} \left(\mathbf{x}^{(k)} - t_k \nabla g(\mathbf{x}^{(k)}) \right).$$

This is also called *forward-backward splitting*, with the “forward” referring to the gradient step, and the “backward” to the proximal step. (The prox step is still making progress, just like the gradient step; the forward and backward refer to the interpretations of gradient descent and the proximal algorithm as forward and backward Euler discretizations, respectively.)

Example: the LASSO

Recall the LASSO, for finding sparse solutions to systems of linear equations:

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \tau \|\mathbf{x}\|_1$$

We take

$$g(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2, \quad \text{so} \quad \nabla g(\mathbf{x}) = \mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{y}),$$

and

$$h(\mathbf{x}) = \tau \|\mathbf{x}\|_1.$$

The prox operator for the ℓ_1 norm is something familiar:

$$\begin{aligned} \text{prox}_{th}(\mathbf{z}) &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(\tau \|\mathbf{x}\|_1 + \frac{1}{2t} \|\mathbf{x} - \mathbf{z}\|_2^2 \right) \\ &= T_{\tau t}(\mathbf{x}), \end{aligned}$$

where $T_{\tau t}$ is the soft-thresholding operator

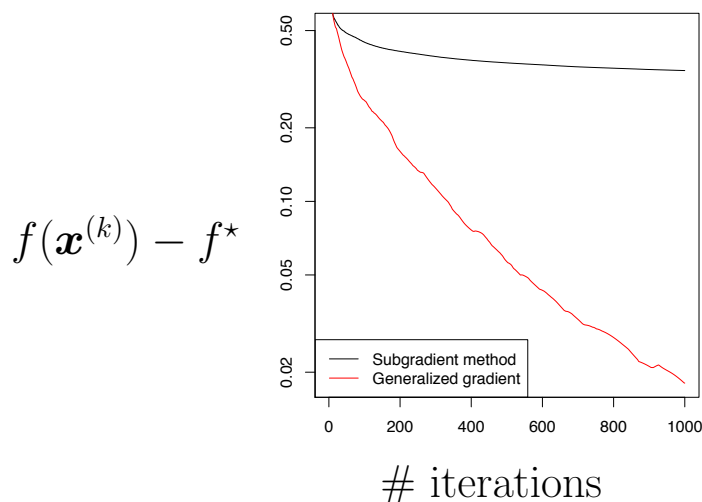
$$(T_\lambda(\mathbf{x})) [i] = \begin{cases} x[i] - \lambda, & x[i] \geq \lambda, \\ 0, & |x[i]| \leq \lambda, \\ x[i] + \lambda, & x[i] \leq -\lambda, \end{cases}$$

with $\lambda = \tau t$. So the gradient step requires an application of \mathbf{A} and \mathbf{A}^T , and the proximal step simply requires a soft-thresholding operation. The iteration looks like

$$\mathbf{x}^{(k+1)} = T_{\tau t_k} \left(\mathbf{x}^{(k)} + t_k \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{x}^{(k)}) \right).$$

This is also called the *iterative soft thresholding algorithm*, or ISTA.

Here is a comparison¹ of a typical run for ISTA versus the subgradient method.



As we can see, ISTA absolutely crushes the subgradient method.

Example: Projected Gradient Descent for Constrained Optimization

Consider the constrained optimization program

$$\underset{\mathbf{x} \in \mathcal{C}}{\text{minimize}} \quad g(\mathbf{x}), \quad (1)$$

where \mathcal{C} is a convex set and g is a differentiable convex function. This program can be re-written as the unconstrained program

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad g(\mathbf{x}) + h(\mathbf{x})$$

where $h(\mathbf{x})$ is the characteristic function for \mathcal{C} :

$$h(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \in \mathcal{C}, \\ +\infty, & \mathbf{x} \notin \mathcal{C}. \end{cases}$$

¹This is taken from the lecture notes of Geoff Gordon and Ryan Tibshirani; “generalized gradient” in the legend means ISTA.

As we have seen before, $h(\mathbf{x})$ is convex. Its prox function is

$$\begin{aligned}\text{prox}_{tf}(\mathbf{z}) &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(h(\mathbf{x}) + \frac{1}{2t} \|\mathbf{x} - \mathbf{z}\|_2^2 \right) \\ &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{x} - \mathbf{z}\|_2^2 \\ &= P_{\mathcal{C}}(\mathbf{z}).\end{aligned}$$

That is, for any $t > 0$, $\text{prox}_{tf}(\mathbf{z})$ is the closest point in \mathcal{C} to \mathbf{z} — it is a projection onto the set \mathcal{C} .

The iteration for solving (1) is then

$$\mathbf{x}^{(k+1)} = P_{\mathcal{C}} \left(\mathbf{x}^{(k)} - t_k \nabla g(\mathbf{x}^{(k)}) \right).$$

That is, we take a gradient step on g , then re-project onto the constraint set. Notice that for $k \geq 1$, the current iterate $\mathbf{x}^{(k)}$ will always be feasible.

This is called **projected gradient descent**. It is a very simple (but effective) method for solving constrained optimization problems when we have an efficient method for projecting onto the constraint set \mathcal{C} .

Example: Least-squares with positivity constraints

Suppose we want to solve

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 \quad \text{subject to} \quad \mathbf{x} \geq \mathbf{0}.$$

The projector onto the set of all positive vectors is simple, we simply just set all of the negative entries to zero. The iteration is

$$\mathbf{x}^{(k+1)} = \left(\mathbf{x}^{(k)} - t_k \mathbf{A}^T (\mathbf{y} - \mathbf{A}\mathbf{x}^{(k)}) \right)_+,$$

where

$$(\mathbf{z})_+[i] = \begin{cases} z[i], & z[i] \geq 0, \\ 0, & z[i] < 0. \end{cases}$$

Example: Nuclear norm minimization (matrix LASSO)

A common problem in machine learning and signal processing is to recover a low-rank matrix from indirect linear measurements. Just as we use the ℓ_1 norm as a proxy for the sparsity of a vector, we can use the *nuclear norm* (sum of singular values) as a proxy for the rank of a matrix.

FINISH THIS ...

Convergence analysis

The convergence analysis is very similar to what we did for the proximal point method. We assume that the gradient of g is Lipschitz, meaning that for any $\mathbf{x} \in \mathbb{R}^N$,

$$g(\mathbf{z}) \leq g(\mathbf{x}) + \nabla g(\mathbf{x})^\top (\mathbf{z} - \mathbf{x}) + \frac{L}{2} \|\mathbf{z} - \mathbf{x}\|_2^2, \quad \text{for all } \mathbf{z} \in \mathbb{R}^N. \quad (2)$$

We will use a fixed “step size”, $t = t_k$ for all k , and we will take $t \leq 1/L$. We will also assume that f has a minimizer \mathbf{x}^* , and we will use the notation $f^* = f(\mathbf{x}^*)$.

Let’s again look at what happens in the iteration moving from $\mathbf{x}^{(k)}$ to $\mathbf{x}^{(k+1)}$. As before, let \mathbf{d} be the direction we move between these two iterations:

$$\mathbf{d} = \frac{1}{t} \left(\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)} \right),$$

so $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - t\mathbf{d}$.

Now

$$\begin{aligned} f(\mathbf{x}^{(k+1)}) &= g(\mathbf{x}^{(k+1)}) + h(\mathbf{x}^{(k+1)}) \\ &\leq g(\mathbf{x}^{(k)}) - t\nabla g(\mathbf{x}^{(k)})^\top \mathbf{d} + \frac{Lt^2}{2} \|\mathbf{d}\|_2^2, \end{aligned} \quad (3)$$

where we have applied the Lipschitz bound (2). We now use two facts to get an upper bound on this expression for arbitrary points \mathbf{z} . First, since $\nabla g(\mathbf{x}^{(k)})$ is a subgradient of g at $\mathbf{x}^{(k)}$,

$$g(\mathbf{z}) \geq g(\mathbf{x}^{(k)}) + \nabla g(\mathbf{x}^{(k)})^\top (\mathbf{z} - \mathbf{x}^{(k)}). \quad (4)$$

Second, since

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \text{prox}_{th}(\mathbf{x}^{(k)} - t\nabla g(\mathbf{x}^{(k)})) \\ &= \arg \min_{\mathbf{x}} \left(h(\mathbf{x}) + \frac{1}{2t} \|\mathbf{x} - \mathbf{x}^{(k)} + t\nabla g(\mathbf{x}^{(k)})\|_2^2 \right), \end{aligned}$$

we know

$$\mathbf{0} \in \partial h(\mathbf{x}^{(k+1)}) - \mathbf{d} + \nabla g(\mathbf{x}^{(k)}) \quad \Rightarrow \quad \mathbf{d} - \nabla g(\mathbf{x}^{(k)}) \in \partial h(\mathbf{x}^{(k+1)}).$$

Thus

$$h(\mathbf{z}) \geq h(\mathbf{x}^{(k+1)}) + (\mathbf{d} - \nabla g(\mathbf{x}^{(k)}))^{\top} (\mathbf{z} - \mathbf{x}^{(k+1)}). \quad (5)$$

We combine (4) and (5) back into (3) to get

$$\begin{aligned} f(\mathbf{x}^{(k+1)}) &\leq g(\mathbf{z}) - \nabla g(\mathbf{x}^{(k)})^{\top} (\mathbf{z} - \mathbf{x}^{(k)}) - t \nabla g(\mathbf{x}^{(k)})^{\top} \mathbf{d} + \frac{Lt^2}{2} \|\mathbf{d}\|_2^2 \\ &\quad + h(\mathbf{z}) - (\mathbf{d} - \nabla g(\mathbf{x}^{(k)}))^{\top} (\mathbf{z} - \mathbf{x}^{(k)} + t\mathbf{d}) \\ &= f(\mathbf{z}) - \mathbf{d}^{\top} (\mathbf{z} - \mathbf{x}^{(k)}) + t \left(\frac{Lt}{2} - 1 \right) \|\mathbf{d}\|_2^2 \\ &\leq f(\mathbf{z}) - \mathbf{d}^{\top} (\mathbf{z} - \mathbf{x}^{(k)}) - \frac{t}{2} \|\mathbf{d}\|_2^2, \end{aligned}$$

where we have used the bound on the step size $t \leq 1/L$ in the last step. Note that the expression above holds for all $\mathbf{z} \in \mathbb{R}^N$.

Evaluating the expression for the upper bound above at $\mathbf{z} = \mathbf{x}^{(k)}$ to see that the functional value is non-increasing at every step:

$$f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}) - \frac{t}{2} \|\mathbf{d}\|_2^2 \leq f(\mathbf{x}^{(k)}).$$

We can also evaluate the upper bound at $\mathbf{z} = \mathbf{x}^*$,

$$f(\mathbf{x}^{(k+1)}) - f^* \leq -\mathbf{d}^{\top} (\mathbf{x}^* - \mathbf{x}^{(k)}) - \frac{t}{2} \|\mathbf{d}\|_2^2,$$

and then complete the square as was done on page 8 of the previous set of notes to yield

$$f(\mathbf{x}^{(k+1)}) - f^* \leq \frac{1}{2t} \left(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_2^2 \right),$$

which again results in the convergence bound

$$\begin{aligned} f(\mathbf{x}^{(k)}) - f^* &\leq \frac{1}{k} \sum_{i=1}^k f(\mathbf{x}^{(i)}) - f^* \\ &\leq \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2}{2kt} \end{aligned}$$

So the proximal gradient algorithm exhibits a convergence rate of $O(1/k)$.

The above is kind of a “master result” for the convergence rate of three different algorithms that are special cases of proximal gradient:

- gradient descent (take $h(\mathbf{x}) = 0$),
- conditional gradient descent (take $h(\mathbf{x})$ as the characteristic function for a set \mathcal{C}),
- the proximal point method (take $g(\mathbf{x}) = 0$).

The work above gives a unified analysis for all three of these, showing that they are all $O(1/k)$.

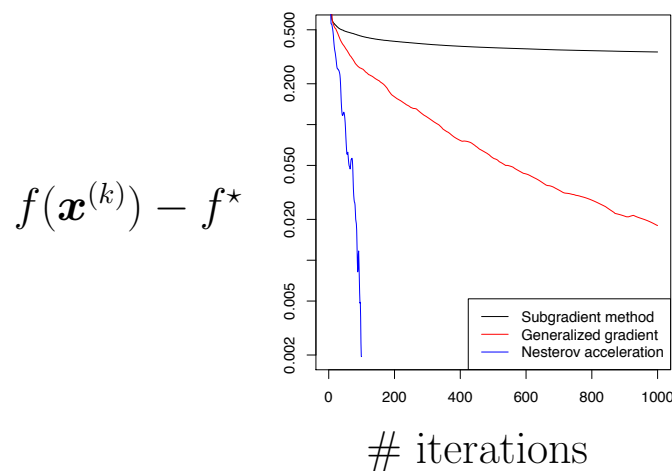
Accelerated proximal gradient

We can accelerate the proximal gradient method in exactly the same way we accelerated the proximal point method — in fact, the ppm acceleration is simply a special case as that for the proximal gradient algorithm. The accelerated iteration is

$$\begin{aligned} \mathbf{w}^{(k)} &= \mathbf{x}^{(k)} + \frac{k-1}{k+2} (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) \\ \mathbf{x}^{(k+1)} &= \text{prox}_{t_k h} (\mathbf{w}^{(k)} - t_k \nabla g(\mathbf{w}^{(k)})). \end{aligned}$$

Again, the actual computations here are no more involved than for the non-accelerated version, but the gains in efficiency (in terms of number of iterations to convergence) can be significant. We will not prove it here (see [BT09] for a complete analysis), but adding in the momentum term above results in convergence rate of $O(1/k^2)$.

The numerical performance can also be dramatically better. Here are typical runs² for the LASSO, which compares the standard proximal gradient method (ISTA) to its accelerated version (FISTA):



References

- [BT09] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.*, 2(1):183–202, 2009.

²Again, this example comes from Gordon and Tibshirani; as before “generalized gradient” means ISTA, and “Nesterov acceleration” means FISTA.