

Proximal Algorithms

The subgradient algorithm is one generalization of gradient descent. It is simple, but the convergence is typically very slow (and it does not even converge in general for a fixed stepsize). The essential reason for this was that there are plenty of subgradients that are large near and even at the solution.

In this section, we will look at another iteration that finds the solution to an unconstrained convex program. The **proximal algorithm** or **proximal point method** can be understood in many different ways, but we will start by thinking about as a regularizer that naturally damps its influence as the iterations proceed, and as a twist on gradient descent.

Our goal is to solve the unconstrained optimization program

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} f(\mathbf{x}) \quad (1)$$

where f is convex but not necessarily smooth. Proximal point methods use the following iteration:

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(f(\mathbf{x}) + \frac{1}{2t_k} \|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2 \right). \quad (2)$$

When f is convex, $f(\mathbf{x}) + \alpha \|\mathbf{x} - \mathbf{z}\|_2^2$ is strictly convex for all $\alpha > 0$ and $\mathbf{z} \in \mathbb{R}^N$, so the mapping from $\mathbf{x}^{(k)}$ to $\mathbf{x}^{(k+1)}$ is well-defined. We will sometimes use the “prox operator” to denote this mapping:

$$\text{prox}_{t f}(\mathbf{z}) = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(f(\mathbf{x}) + \frac{1}{2t} \|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2 \right).$$

At this point, you would be forgiven for wondering what all the fuss is about. We have taken the unconstrained problem (1) and turned

it into a series of constrained problems (2). For this to make sense, this program would have to be easier to solve for some reason. This can certainly be the case, and we will see an example in the next section. One way to think about the additional $\frac{1}{2t_k} \|\mathbf{x} - \mathbf{x}^{(k)}\|$ is as a *regularizer* whose influence naturally disappears as we approach the solution, even for a fixed “step size” $t_k = t$. Computationally, the smoothed problem can be much easier to solve.

But past this, we will use our discussion here to build up to the so-called *proximal gradient* algorithm, discussed in the next set of notes, which is a great way to minimize an unconstrained function that can be broken into a smooth part and a nonsmooth part. A very nice detailed review of proximal algorithms can be found in [PB14].

Least squares

Suppose we want to solve the standard least-squares problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2.$$

When \mathbf{A} has full column rank, we know that the solution is given by $\hat{\mathbf{x}}_{\text{ls}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$. However, we also know that when $\mathbf{A}^T \mathbf{A}$ is not well-conditioned, this inverse can be unstable to compute, and iterative descent methods (gradient descent and conjugate gradients) can take many iterations to converge.

Consider the proximal point iteration (with fixed $t_k = t$) for solving this problem:

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(\frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \frac{1}{2t} \|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2 \right).$$

Here we have the closed form solution

$$\begin{aligned}\mathbf{x}^{(k+1)} &= (\mathbf{A}^T \mathbf{A} + \delta \mathbf{I})^{-1} (\mathbf{A}^T \mathbf{y} + \delta \mathbf{x}^{(k)}), \quad \delta = \frac{1}{t} \\ &= \mathbf{x}^{(k)} + (\mathbf{A}^T \mathbf{A} + \delta \mathbf{I})^{-1} \mathbf{A}^T (\mathbf{y} - \mathbf{A} \mathbf{x}^{(k)}).\end{aligned}$$

Now each step is equivalent to solving a least-squares problem, but this problem can be made well-conditioned by choosing δ (t) appropriately. The iterations above will converge to $\hat{\mathbf{x}}_{\text{ls}}$ for any value of t ; as we decrease t (increase δ), the number of iterations to get within a certain accuracy of $\hat{\mathbf{x}}_{\text{ls}}$ increases, but the least-squares problems involved are all very well conditioned. For t very small, we are back at gradient descent (with stepsize t).

This is actually a well-known technique in numerical linear algebra called *iterative refinement*.

Implicit gradient descent (“backward Euler”)

The proximal point method can also be interpreted as a variation on gradient descent. To see this, let us return for a moment to the differential equations for the “gradient flow” of f :

$$\mathbf{x}'(t) = -\nabla f(\mathbf{x}), \quad \mathbf{x}(0) = \mathbf{x}_0. \quad (3)$$

The equilibrium points for this system are the \mathbf{x} such that $\nabla f(\mathbf{x}) = \mathbf{0}$, which are precisely the minimizers for $f(\mathbf{x})$.

As we first saw in the notes on the heavy ball method, we can interpret gradient descent as a first-order numerical method for tracing the path from \mathbf{x}_0 to a solution \mathbf{x}^* . This comes from discretizing the derivative on the right using a forward finite difference:

$$\frac{\mathbf{x}(t+h) - \mathbf{x}(t)}{h} \approx -\nabla f(\mathbf{x}(t)) \quad \text{for small } h.$$

Thus the gradient descent iterations

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - h \nabla f(\mathbf{x}^{(k)})$$

approximate the solution at equispaced times spaced h seconds apart — the stepsize in gradient descent can be interpreted as the time scale to which we are approximating the derivative. This is known as the *forward Euler method* for discretizing (3).

But now suppose we used a *backward difference* to approximate the derivative:

$$\frac{\mathbf{x}(t) - \mathbf{x}(t - h)}{h} \approx -\nabla f(\mathbf{x}(t)) \quad \text{for small } h.$$

Now the iterates must obey

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - h \nabla f(\mathbf{x}^{(k+1)}).$$

This is known as the *backward Euler method* for discretizing (3). Computing the iterates is not as straightforward — we can just compute the gradient at the current point, we have to find the next point by finding a $\mathbf{x}^{(k+1)}$ that obeys the equation above.

This is exactly what the proximal operator does. If f is differentiable, then

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(f(\mathbf{x}) + \frac{1}{2t} \|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2 \right) \\ &\quad \Downarrow \\ \mathbf{0} &= \nabla f(\mathbf{x}^{(k+1)}) + \frac{1}{t} (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}). \end{aligned} \quad (4)$$

So the proximal point method can be interpreted as a backward Euler discretization for gradient flow.

Note that we assumed the differentiability of f above purely for illustration; we can compute the prox operator whether or not f has a gradient.

Convergence

We have yet to show that the iteration

$$\mathbf{x}^{(k+1)} = \text{prox}_{tf}(\mathbf{x}^{(k)})$$

actually converges to a minimizer of f . We will do this now, and see that the convergence guarantees (at least in terms of number of iterations) are much nicer than they are for the subgradient method.

First question: If we are at a solution \mathbf{x}^* , does applying the proximal operator move us off of this solution?

Answer: No. \mathbf{x}^* is a minimizer of convex $f(\mathbf{x})$ if and only if

$$\mathbf{x}^* = \text{prox}_{tf}(\mathbf{x}^*)$$

for every $t > 0$. That is, minimizers of f are fixed points of the proximal operator. Just like applying another iteration of gradient descent at the minimizer of a smooth function moves you nowhere (since the gradient is zero), applying the prox operator to the minimizer of a general convex function does nothing. This is in stark contrast to the subgradient method, where we can definitely move off a solution.

The fixed-point property can be established with a quick argument. If \mathbf{x}^* is a minimizer of f , then

$$\begin{aligned} \mathbf{0} &\in \partial f(\mathbf{x}^*) \\ &= \partial f(\mathbf{x}^*) + \frac{1}{t}(\mathbf{x}^* - \mathbf{x}^*) \end{aligned}$$

and so $\mathbf{0}$ is also in the subdifferential of $f(\mathbf{x}) + \frac{1}{2t}\|\mathbf{x} - \mathbf{x}^*\|_2^2$ at $\mathbf{x} = \mathbf{x}^*$. The converse is now also clear, as if $\mathbf{0}$ is in the subdifferential of $f(\mathbf{x}) + \frac{1}{2t}\|\mathbf{x} - \mathbf{x}^*\|_2^2$ at $\mathbf{x} = \mathbf{x}^*$, then

$$\begin{aligned} \mathbf{0} &\in \partial f(\mathbf{x}^*) + \frac{1}{t}(\mathbf{x}^* - \mathbf{x}^*) \\ &\Downarrow \\ \mathbf{0} &\in \partial f(\mathbf{x}^*), \end{aligned}$$

and so \mathbf{x}^* is a minimizer of f .

We can also bound the number of steps it takes for the proximal point method to achieve a certain accuracy, just as we have for the other unconstrained minimization algorithms that we have encountered. Under the conditions that f is a closed convex function (so that the prox function is well-defined), and that there exists a minimizer \mathbf{x}^* , we will show below that the iterations $\mathbf{x}^{(k+1)} = \text{prox}_{t_k f}(\mathbf{x}^{(k)})$ obey

$$f(\mathbf{x}^{(k)}) - f^* \leq \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2}{2 \sum_{i=1}^k t_i} \quad \text{for all } k \geq 1, \quad (5)$$

where $f^* = f(\mathbf{x}^*)$. From here we can see that

$$t_k = t \quad \Rightarrow \quad f(\mathbf{x}^{(k)}) - f^* = O(1/k),$$

that is, we are guaranteed to get to a point with $\mathbf{x}^{(k)}$ such that $f(\mathbf{x}^{(k)}) - f^* \leq \epsilon$ in $O(1/\epsilon)$ iterations.

In looking at (5), it seems that we can make the convergence as fast as we want by making the t_k very large. This is true, but remember that each step itself involves solving an optimization program, and the cost of solving this program might rely critically on the t_k . As t_k gets large, we are effectively solving the original program itself.

We will establish (5) in three steps:

1. We will show that $f(\mathbf{x}^{(i)})$ is non-increasing
2. We will also show that

$$f(\mathbf{x}^{(i)}) - f^* \leq \frac{1}{2t_i} \left(\|\mathbf{x}^{(i-1)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 \right)$$

3. Combining the two facts above gives us

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

We will start by looking at what happens at a single iteration. Define

$$\mathbf{g}^{(i)} = \frac{1}{t} \left(\mathbf{x}^{(i)} - \text{prox}_{tf}(\mathbf{x}^{(i)}) \right).$$

Thus $\mathbf{g}^{(i)}$ is the (negative of) the direction we move in going from $\mathbf{x}^{(i)}$ to $\mathbf{x}^{(i+1)}$:

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - t\mathbf{g}^{(i)}.$$

Since we are just going to study this one iteration, we will make the notation cleaner by dropping the superscript, taking $\mathbf{g} := \mathbf{g}^{(i)}$ and $t = t_i$ below.

Since

$$\mathbf{x}^{(i+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(f(\mathbf{x}) + \frac{1}{2t} \|\mathbf{x} - \mathbf{x}^{(i)}\|_2^2 \right),$$

the optimality conditions for the program on the right tell us that

$$\mathbf{0} \in \partial f(\mathbf{x}^{(i+1)}) + \frac{1}{t} \left(\mathbf{x}^{(i+1)} - \mathbf{x}^{(i)} \right),$$

which means

$$\mathbf{g} \in \partial f(\mathbf{x}^{(i+1)}).$$

By the definition of subgradient, then,

$$\begin{aligned} f(\mathbf{x}^{(i+1)}) &\leq f(\mathbf{z}) + \mathbf{g}^\top(\mathbf{x}^{(i+1)} - \mathbf{z}) \\ &= f(\mathbf{z}) + \mathbf{g}^\top(\mathbf{x}^{(i)} - \mathbf{z} - t\mathbf{g}) \\ &= f(\mathbf{z}) + \mathbf{g}^\top(\mathbf{x}^{(i)} - \mathbf{z}) - t\|\mathbf{g}\|_2^2 \quad \text{for all } \mathbf{z} \in \mathbb{R}^N. \end{aligned} \quad (6)$$

In particular, if we evaluate the above at $\mathbf{x} = \mathbf{x}^{(i)}$, then

$$f(\mathbf{x}^{(i+1)}) \leq f(\mathbf{x}^{(i)}) - t\|\mathbf{g}\|_2^2.$$

Thus $f(\mathbf{x}^{(i+1)}) \leq f(\mathbf{x}^{(i)})$, and property 1 above holds.

To get property 2, we evaluate (6) at $\mathbf{z} = \mathbf{x}^*$ to get

$$f(\mathbf{x}^{(i+1)}) - f^* \leq \mathbf{g}^\top(\mathbf{x}^{(i)} - \mathbf{x}^*) - t\|\mathbf{g}\|_2^2.$$

Now we “complete the square” on the right. Since $t > 0$,

$$\begin{aligned} f(\mathbf{x}^{(i+1)}) - f^* &\leq \mathbf{g}^\top(\mathbf{x}^{(i)} - \mathbf{x}^*) - \frac{t}{2}\|\mathbf{g}\|_2^2 \\ &= \frac{1}{2t} \left(2t\mathbf{g}^\top(\mathbf{x}^{(i)} - \mathbf{x}^*) - t^2\|\mathbf{g}\|_2^2 \right) \\ &= \frac{1}{2t} \left(-\|\mathbf{x}^{(i)} - \mathbf{x}^* - t\mathbf{g}\|_2^2 + \|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 \right) \\ &= \frac{1}{2t} \left(\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2 \right). \end{aligned}$$

Thus property 2 above holds.

Accelerated proximal points algorithms

We can “accelerate” the proximal point method in the same way we accelerated gradient descent: by adding “feedback” from previous iterations. This result is essentially contained in [BT09]; we will discuss it more in the next section when we talk about the proximal gradient algorithm.

The essential iteration is as follows:

$$\begin{aligned}\mathbf{x}^{(1)} &= \text{prox}_{t_1 f}(\mathbf{x}^{(0)}), \\ \mathbf{x}^{(k)} &= \text{prox}_{t_k f} \left(\mathbf{x}^{(k-1)} + s_k(\mathbf{x}^{(k-1)} - \mathbf{x}^{(k-2)}) \right),\end{aligned}$$

where the t_k are given as before, and the s_k are computed as

$$s_k = \theta_k \left(\frac{1}{\theta_{k-1}} - 1 \right), \quad \text{where } \theta_k \text{ obeys } \frac{\theta_k^2}{\theta_{k-1}^2} = (1 - \theta_k) \frac{t_k}{t_{k-1}}.$$

So we are adding a little bit of momentum to $\mathbf{x}^{(k-1)}$ before we compute the prox mapping. This makes sense intuitively, and we have seen that this type of idea leads to real gains in standard gradient descent. Deriving the step sizes above, though, is non-trivial, and we will omit it here.

The convergence results for the accelerated proximal method say that the iteration converges if the “step sizes” t_k are chosen such that $\sum_i \sqrt{t_i} \rightarrow \infty$. For constant step size, we have $O(1/k^2)$ convergence, which is a significant improvement over the standard proximal point method.

Moreover, this acceleration comes at practically zero additional computational cost. All we have to do is store an extra iterate in memory.

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.
- [PB14] N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends in Optimization*, 1(3):123–231, 2014.