

## 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 step size).

One way to deal with this is to add a smooth *regularization* term. Specifically, it is easy to show that if  $\mathbf{x}^*$  is a minimizer of  $f(\mathbf{x})$ , then it is also the minimizer of

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}) + \delta \|\mathbf{x} - \mathbf{x}^*\|_2^2,$$

where  $\delta > 0$ . While the resulting optimization problem is still nonsmooth, it is now strongly convex, and we know that strongly convex functions are generally much easier to minimize. The “only” challenge is that it requires us to already know the solution  $\mathbf{x}^*$ , which would seem to limit the practical applicability of this idea.

We can turn this into an actual algorithm by adopting an iterative approach. The **proximal algorithm** or **proximal point method** uses the following iteration:<sup>1</sup>

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

As noted above, when  $f$  is convex,  $f(\mathbf{x}) + \delta \|\mathbf{x} - \mathbf{z}\|_2^2$  is strictly convex for all  $\delta > 0$  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:

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

---

<sup>1</sup>The notation  $\arg \min_{\mathbf{x} \in \mathbb{R}^N} f(\mathbf{x})$  denotes a function that returns the vector  $\mathbf{x} \in \mathbb{R}^N$  that minimizes  $f(\mathbf{x})$ . This is in contrast to  $\min_{\mathbf{x} \in \mathbb{R}^N} f(\mathbf{x})$ , which refers to the minimum possible value for  $f(\mathbf{x})$ .

At this point, you would be forgiven for having doubts about what we are really doing here. We have taken an optimization problem and turned it into... a sequence of *many* optimization problems. However, these problems can sometimes be far easier to solve than the original problem. One way to think about the additional  $\frac{1}{2\alpha_k}\|\mathbf{x} - \mathbf{x}_k\|$  term is as a *regularizer* that makes each subproblem computationally easier to solve, and whose influence naturally disappears as we approach the solution, even for a fixed “step size”  $\alpha_k = \alpha$ .

## 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  $f$ :

$$\mathbf{x}'(t) = -\nabla f(\mathbf{x}(t)), \quad \mathbf{x}(0) = \mathbf{x}_0. \quad (2)$$

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 discussed in the context of momentum-based methods, 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 step size 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 (2).

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 an equally valid technique for discretizing (2) known as the *backward Euler method*. However, computing the iterates is not as straightforward – we can't just compute the gradient at the current point, we have to find the next point by finding an  $\mathbf{x}_{k+1}$  that obeys the equation above.

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

$$\begin{aligned} \mathbf{x}_{k+1} &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left( f(\mathbf{x}) + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right) \\ &\quad \Downarrow \\ \mathbf{0} &= \nabla f(\mathbf{x}_{k+1}) + \frac{1}{\alpha_k} (\mathbf{x}_{k+1} - \mathbf{x}_k). \end{aligned} \tag{3}$$

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.

## Example: Least squares

Suppose we want to solve the standard least-squares problem

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

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  $\alpha_k = \alpha$ ) for solving this problem:

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

Here we have the closed form solution

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

Now each step is equivalent to solving a least-squares problem, but this problem can be made well-conditioned by choosing  $\delta$  (i.e.,  $\alpha$ ) appropriately. The iterations above will converge to  $\hat{\mathbf{x}}_{\text{ls}}$  for any value of  $\alpha$ ; as we decrease  $\alpha$  (increase  $\delta$ ), 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  $\alpha$  very small, we are back at gradient descent (with step size  $\alpha$ ).

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

## Proximal gradient algorithms

Recall the core update equation for the proximal point method:

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

Suppose that we did not wish to fully solve this problem at each iteration. If  $f$  is differentiable, we could approximate this update by replacing  $f(\mathbf{x})$  with its linear approximation  $f(\mathbf{x}_k) + \langle \mathbf{x} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle$ . This would yield the update

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

Thus, taking a linear approximation of  $f$ , the proximal method simply reduces to standard gradient descent.

Where this starts getting interesting is when we encounter optimization problems where the objective function can be broken into the sum two parts, i.e.,

$$f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x}),$$

where both  $g$  and  $h$  are convex, but  $g$  is smooth (differentiable) and  $h$  is a non-smooth function for which there is a fast proximal operator. Such optimization problems quite a bit more often than you might expect.

The **proximal gradient** algorithm is the result of applying the proximal point method to minimize the approximation of  $f$  where we take a linear approximation to the smooth component  $g$ . Using the same argument as above, this results in the update rule

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

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 our friend the LASSO:

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

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  $\ell_1$  norm is:

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

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

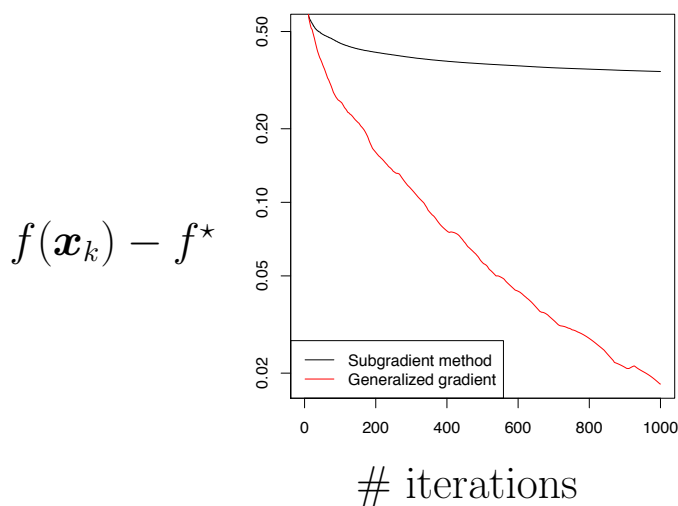
$$[T_{\tau\alpha}(\mathbf{z})]_n = \begin{cases} z_n - \tau\alpha, & z_n \geq \tau\alpha, \\ 0, & |z_n| \leq \tau\alpha, \\ z_n + \tau\alpha, & z_n \leq -\tau\alpha. \end{cases}$$

Hence, 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

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

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

Here is a comparison<sup>2</sup> of a typical run for ISTA versus the subgradient method. ISTA absolutely crushes the subgradient method.



<sup>2</sup>This is taken from the lecture notes of Geoff Gordon and Ryan Tibshirani; “generalized gradient” in the legend means ISTA.

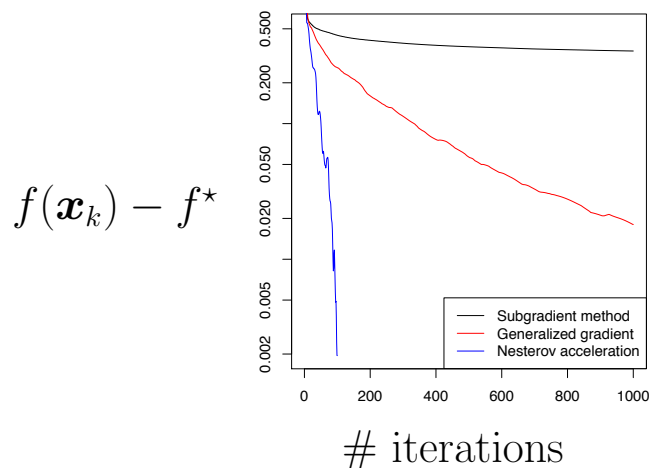
## Accelerated proximal gradient

We can accelerate the proximal gradient method in exactly the same way we accelerated gradient descent – in fact, the Nesterov’s method for gradient descent is simply a special case as that for the proximal gradient algorithm. The accelerated iteration is

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

Again, the computations here are in general no more involved than for the non-accelerated version, but the number of iterations can be significantly lower. We will not prove it here (the proof is straightforward, but long), but adding in the momentum term results in convergence rate of  $O(1/k^2)$  using a similar argument as before.

The numerical performance can also be dramatically better. Here are typical runs<sup>3</sup> for the LASSO, which compares the standard proximal gradient method (ISTA) to its accelerated version (FISTA):



<sup>3</sup>Again, this example comes from Gordon and Tibshirani; as before “generalized gradient” means ISTA, and “Nesterov acceleration” means FISTA.



## Convergence of the proximal gradient method

The convergence analysis of the proximal gradient method is extremely similar to what we did for gradient descent. In fact, gradient descent is a special case of the proximal gradient method (when  $h(\mathbf{x}) = 0$ ), and our analysis will recover the same result. We will assume that  $g$  is  $M$ -smooth, but we will make no assumptions on  $h$  aside from convexity. As before, we will use a fixed “step size”,  $\alpha_k = 1/M$  for all  $k$ . We will  $\mathbf{x}^*$  denote any minimizer of  $f$ .

The general structure of the argument is as follows:

1. Using the  $M$ -smoothness of  $g$  as well as the first-order characterization of convexity, we can establish that

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{z}) + \langle \mathbf{x}_k - \mathbf{z}, \mathbf{d}_k \rangle - \frac{1}{2M} \|\mathbf{d}_k\|_2^2 \quad (4)$$

for all  $\mathbf{z} \in \mathbb{R}^N$  where  $\mathbf{d}_k := M(\mathbf{x}_k - \mathbf{x}_{k+1})$ .

2. From (4) we can conclude, by setting  $\mathbf{z} = \mathbf{x}_k$ , that

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{1}{2M} \|\mathbf{d}_k\|_2^2 \leq f(\mathbf{x}_k),$$

and thus  $f(\mathbf{x}_k)$  is non-increasing at every step.

3. From (4) we can also conclude, by setting  $\mathbf{z} = \mathbf{x}^*$ , that

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}^*) + \langle \mathbf{x}_k - \mathbf{x}^*, \mathbf{d}_k \rangle - \frac{1}{2M} \|\mathbf{d}_k\|_2^2.$$

By exactly the same argument as we have seen in the analysis of both gradient descent and Nesterov’s method, we can show that this bound is equivalent to

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \frac{M}{2} (\|\mathbf{x}_k - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2).$$

4. This yields a telescopic sum, and hence by an identical argument to that used in analyzing gradient descent, we arrive at the bound

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{M}{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$

Thus, the proximal gradient algorithm exhibits the same convergence rate as gradient descent:  $O(1/k)$ . This is remarkable when considering that it holds for *any*  $h$ . This result is in fact a kind of “master result” for the convergence rate of many different algorithms:

- gradient descent (take  $h(\mathbf{x}) = 0$ ),
- the proximal point method (take  $g(\mathbf{x}) = 0$ ),
- the proximal gradient method.

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

Note that the only novelty in the analysis above compared to that of gradient descent is the derivation of (4). To establish this inequality, we first note that

$$\begin{aligned} f(\mathbf{x}_{k+1}) &= g(\mathbf{x}_{k+1}) + h(\mathbf{x}_{k+1}) \\ &\leq g(\mathbf{x}_k) - \frac{1}{M} \langle \mathbf{d}_k, \nabla g(\mathbf{x}_k) \rangle + \frac{1}{2M} \|\mathbf{d}_k\|_2^2 + h(\mathbf{x}_{k+1}), \end{aligned} \quad (5)$$

where the inequality follows directly from the definition of  $M$ -smoothness. We now use two facts to get an upper bound on this expression. First, note that from the first-order characterization of convexity,

$$g(\mathbf{z}) \geq g(\mathbf{x}_k) + \langle \mathbf{z} - \mathbf{x}_k, \nabla g(\mathbf{x}_k) \rangle. \quad (6)$$

Second, since

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

we know

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

Thus

$$h(\mathbf{z}) \geq h(\mathbf{x}_{k+1}) + \langle \mathbf{z} - \mathbf{x}_{k+1}, \mathbf{d}_k - \nabla g(\mathbf{x}_k) \rangle. \quad (7)$$

We combine (6) and (7) back into (5) to obtain

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

which establishes (4).