

## Accelerated first-order methods

We have seen three techniques for solving an unconstrained, smooth optimization program: gradient descent, Newton’s method, and BFGS, which was an example of a quasi-Newton method. Gradient descent had the advantage of only needing to compute “first order” derivatives of the function we are minimizing (e.g. the gradient) — BFGS also has this property, but was implicitly trying to use changes in the first derivative to derive a second-order model.

There are small changes we can make to gradient descent that can dramatically improve its performance, both in theory and in practice. We talk about two of these here: the heavy ball method, and Nesterov’s “optimal algorithm”.

### The Heavy Ball method

One way to interpret gradient descent is as a discretization to the *gradient flow* differential equation

$$\begin{aligned}\mathbf{x}'(t) &= -\nabla f(\mathbf{x}(t)), \\ \mathbf{x}(0) &= \mathbf{x}_0.\end{aligned}$$

The solution to these equations is a curve that tracks the direction of steepest descent directly to the minimizer. To see how gradient descent arises, we discretize the derivative with a forward difference

$$\mathbf{x}'(t) \approx \frac{\mathbf{x}(t+h) - \mathbf{x}}{h},$$

for some small  $h$ . So if we think of  $\mathbf{x}^{(k+1)}$  and  $\mathbf{x}^{(k)}$  as closely spaced time points, we can interpret

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

as a discrete approximation to gradient flow. Re-arranging the equation above yields the gradient descent iteration  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)})$ .

The problem is once we perform this discretization, the path tends to oscillate. One way to get a more regular path is to add a second-order term to the differential equation:

$$\mathbf{x}''(t) + a\mathbf{x}'(t) = -b\nabla f(\mathbf{x}(t))$$

From a physical perspective, this is a model for adding friction to a particle moving in potential field. Discretizing the dynamics as before with

$$\mathbf{x}''(t) \approx \frac{\mathbf{x}^{(k+1)} - 2\mathbf{x}^{(k)} + \mathbf{x}^{(k-1)}}{h_1}, \quad \mathbf{x}'(t) \approx \frac{\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}}{h_2}$$

gives us the iteration for the **heavy ball method**, introduced by Polyak in 1964 [Pol64]:

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

The second term above adds a little bit of the last step  $\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}$  direction into the new step direction  $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$  — this method is also referred to as *gradient descent with momentum*.

## Convergence for strongly convex functions

We have seen that gradient descent exhibits linear convergence when  $f(\mathbf{x})$  is strongly convex:

$$m\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq M\mathbf{I}, \quad \text{for all } \mathbf{x} \in \text{dom } f.$$

This means that the eigenvalues of the Hessian matrix are uniformly bounded between  $m$  and  $M$  at all points  $\mathbf{x}$ .

The analysis for gradient descent uses strong convexity with an application of the Taylor theorem. Briefly, we have

$$\begin{aligned}
\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_2 &= \|\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)}) - \mathbf{x}^*\|_2 \\
&= \left\| \mathbf{x}^{(k)} - \mathbf{x}^* - \alpha_k \left( \nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^*) \right) \right\|_2 \\
&= \left\| (\mathbf{I} - \alpha_k \nabla^2 f(\mathbf{z})) (\mathbf{x}^{(k)} - \mathbf{x}^*) \right\|_2 \\
&\leq \|\mathbf{I} - \alpha_k \nabla^2 f(\mathbf{z})\| \cdot \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2,
\end{aligned}$$

where the second equality comes from the fact that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ , while the third equality comes from the Taylor theorem: there exists some  $\mathbf{z}$  on the line between  $\mathbf{x}^{(k)}$  and  $\mathbf{x}^*$  such that

$$\nabla^2 f(\mathbf{z})(\mathbf{x}^{(k)} - \mathbf{x}^*) = \nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^*).$$

Since we have a bound on the eigenvalues of  $\nabla^2 f(\mathbf{z})$ , we know that the maximum eigenvalue of the symmetric matrix  $\mathbf{I} - \alpha_k \nabla^2 f(\mathbf{z})$  is no more than

$$\|\mathbf{I} - \alpha_k \nabla^2 f(\mathbf{z})\| \leq \max(|1 - \alpha_k m|, |1 - \alpha_k M|)$$

If we take  $\alpha_k = 2/(M + m)$ , we have

$$\|\mathbf{I} - \alpha_k \nabla^2 f(\mathbf{z})\| \leq \frac{M - m}{M + m} = \frac{\kappa - 1}{\kappa + 1},$$

where  $\kappa = M/m$  is the “condition number”. So we have

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

which, by induction on  $k$  means

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2 \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^k \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2.$$

When  $\kappa$  is even moderately large, we have

$$\frac{\kappa - 1}{\kappa + 1} \approx 1 - \frac{1}{\kappa},$$

which means that in order to guarantee that we have reduced the initial error by a factor of  $\epsilon$ ,  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2 / \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq \epsilon$ , we need

$$k \gtrsim \kappa \cdot \log(1/\epsilon)$$

For the heavy ball method, we have a similar analysis that ends in a better result. We start by looking at how  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 + \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2$  goes to zero for fixed values of  $\alpha, \beta$  which we will choose later:

$$\begin{aligned} \left\| \begin{bmatrix} \mathbf{x}^{(k+1)} - \mathbf{x}^* \\ \mathbf{x}^{(k)} - \mathbf{x}^* \end{bmatrix} \right\|_2 &= \left\| \begin{bmatrix} \mathbf{x}^{(k)} + \beta(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) - \mathbf{x}^* \\ \mathbf{x}^{(k)} - \mathbf{x}^* \end{bmatrix} - \alpha \begin{bmatrix} \nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^*) \\ \mathbf{0} \end{bmatrix} \right\|_2 \\ &= \left\| \begin{bmatrix} \mathbf{x}^{(k)} + \beta(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) - \mathbf{x}^* \\ \mathbf{x}^{(k)} - \mathbf{x}^* \end{bmatrix} - \alpha \begin{bmatrix} \nabla^2 f(\mathbf{z}) (\mathbf{x}^{(k)} - \mathbf{x}^*) \\ \mathbf{0} \end{bmatrix} \right\|_2 \\ &= \left\| \begin{bmatrix} (1 + \beta)\mathbf{I} - \alpha \nabla^2 f(\mathbf{z}) & -\beta\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}^{(k)} - \mathbf{x}^* \\ \mathbf{x}^{(k-1)} - \mathbf{x}^* \end{bmatrix} \right\|_2. \end{aligned}$$

Applying the above iteratively means that, in the limit

$$\left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|_2 \lesssim \|\mathbf{T}^k\| \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2.$$

where

$$\mathbf{T} = \begin{bmatrix} (1 + \beta)\mathbf{I} - \alpha \nabla^2 f(\mathbf{z}) & -\beta\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}.$$

It is a fundamental result from numerical linear algebra that as  $k$  gets large,  $\|\mathbf{T}^k\|$  becomes very close to<sup>1</sup>  $\rho(\mathbf{T})^k$ , where  $\rho$  is the maximum of the magnitudes of the eigenvalues of  $\mathbf{T}$ .

We can get at these eigenvalues in a systematic way. Start by taking an eigenvalue decomposition of the Hessian matrix above,  $\nabla^2 f(\mathbf{z}) = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$ . Since  $\mathbf{V}\mathbf{V}^T = \mathbf{I}$ , we can write

$$\begin{bmatrix} (1 + \beta)\mathbf{I} - \alpha\nabla^2 f(\mathbf{z}) & -\beta\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{V} & \mathbf{0} \\ \mathbf{0} & \mathbf{V} \end{bmatrix} \begin{bmatrix} (1 + \beta)\mathbf{I} - \alpha\mathbf{\Lambda} & -\beta\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{V}^T \end{bmatrix}.$$

Since  $\begin{bmatrix} \mathbf{V} & \mathbf{0} \\ \mathbf{0} & \mathbf{V} \end{bmatrix}$  is orthonormal, its application on the left and the right of a matrix does not change the magnitude of its eigenvalues, and we have

$$\begin{aligned} \rho(\mathbf{T}) &= \rho\left(\begin{bmatrix} (1 + \beta)\mathbf{I} - \alpha\mathbf{\Lambda} & -\beta\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}\right) \\ &= \max_{n=1,\dots,N} \rho(\mathbf{T}_n), \end{aligned}$$

where

$$\mathbf{T}_n = \begin{bmatrix} 1 + \beta - \alpha\lambda_n & -\beta \\ 1 & 0 \end{bmatrix}.$$

This last equality follows from the fact the large  $2N \times 2N$  matrix can have its rows and columns permuted (which doesn't change the eigenvalues) to become block diagonal, with the  $2 \times 2$  matrices  $\mathbf{T}_n$  as the blocks.

We now have the problem of finding  $\alpha, \beta$  that minimize the size of the largest eigenvalue of the  $2 \times 2$  matrices above given the knowledge

---

<sup>1</sup>If  $\mathbf{T}$  is symmetric, then of course  $\|\mathbf{T}^k\| = \rho(\mathbf{T})^k$  for all  $k$ .

that  $m \leq \lambda_n \leq M$ . A very technical calculation yields the fact that taking

$$\alpha = \frac{4}{(\sqrt{M} - \sqrt{m})^2}, \quad \beta = \frac{\sqrt{M} - \sqrt{m}}{\sqrt{M} + \sqrt{m}},$$

yields

$$\rho(\mathbf{T}_n) \leq \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}, \quad \text{for all } n.$$

Thus the convergence of the heavy ball method is approximately

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2 \lesssim \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2$$

Again, for only moderately large  $\kappa$ , we have

$$\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \approx 1 - \frac{1}{\sqrt{\kappa}},$$

meaning that

$$\frac{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2}{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2} \leq \epsilon \quad \text{when} \quad k \gtrsim \sqrt{\kappa} \log(1/\epsilon).$$

The difference with gradient descent can be significant. When  $\kappa = 10^2$ , we are asking for  $\approx 100 \log(1/\epsilon)$  iterations for gradient descent, as compared with  $\approx 10 \log(1/\epsilon)$  from the heavy ball method.

## References

- [Pol64] B. Polyak. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.