

Optimization using Gradient Descent

Many of the statistical estimation problems that we have talked about ultimately come down to solving an optimization program,

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

for some $f : \mathbb{R}^N \rightarrow \mathbb{R}$. One approach to solving (1) is by using **gradient descent**. The essential algorithm, in the unconstrained case, is straightforward:

From initial point \mathbf{x}_0 , we iterate over $k = 0, 1, \dots$ with

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

The α_k above is a stepsize that changes (in general) at every iteration. When f is *convex*, there is a lot we can say about the convergence of gradient descent to a global optima. When f is not convex, there is much less we can say; most results in this domain characterize the convergence to a local minima that depends on the starting point \mathbf{x}_0 .

Choosing the stepsize

Many algorithms for optimization rely on choosing α_k adaptively. We are at starting point \mathbf{x}_k and have decided to move in direction \mathbf{d}_k , and now we want to decide how far to move; we want to choose α_k for the update $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$. The first thing we might think of is an **exact line search**, taking

$$\alpha_k = \arg \min_{t > 0} f(\mathbf{x}_k + t \mathbf{d}_k).$$

And in fact, this is what we did previously in the least-squares case, as we had a closed form to compute the minimizer exactly. But in general, computing the solution to this 1D optimization program is not

straightforward. Moreover, finding the truly “optimal” step might not even be that beneficial, as \mathbf{d}_k does not point exactly towards the solutions anyway.

More commonly, a step size is chosen through a **backtracking line search**. For a fixed initial step size t_0 and backtracking parameter $0 < \beta < 1$, do the following:

```

Initialize:  $t = t_0$ 
while  $f(\mathbf{x}_k + t\mathbf{d}_k)$  is “not good enough” do
     $t = \beta t$ 
end while
 $\alpha_k = t$ 

```

Of course, the key thing here is defining “not good enough”. You might be tempted to try to implement a condition such as $f(\mathbf{x}_k + t\mathbf{d}_k)$ being as least 10% smaller than $f(\mathbf{x}_k)$, but there is nothing that says that this has to be true for any t , especially when you are close to the solution. When $f(\cdot)$ is convex and the problem is unconstrained, there is a systematic criteria you can use that also gives great results, both in theory and in practice.

For a convex function, it will always be true that

$$f(\mathbf{x}_k + t\mathbf{d}_k) \geq f(\mathbf{x}_k) + t\langle \mathbf{d}_k, \nabla f(\mathbf{x}_k) \rangle, \quad \text{for all } t \geq 0.$$

A good criteria for backtracking line search is to find a t such that

$$f(\mathbf{x}_k + t\mathbf{d}_k) < f(\mathbf{x}_k) + t\gamma\langle \mathbf{d}_k, \nabla f(\mathbf{x}_k) \rangle,$$

for some fixed $\gamma \in (0, 0.5)$. So in the algorithm above, make the substitution

$$f(\mathbf{x}_k + t\mathbf{d}_k) \text{ is “not good enough”} \leftarrow f(\mathbf{x}_k + t\mathbf{d}_k) > f(\mathbf{x}_k) + t\gamma\langle \mathbf{d}_k, \nabla f(\mathbf{x}_k) \rangle.$$

Here is a figure¹ (from [BV04, Chapter 9]):

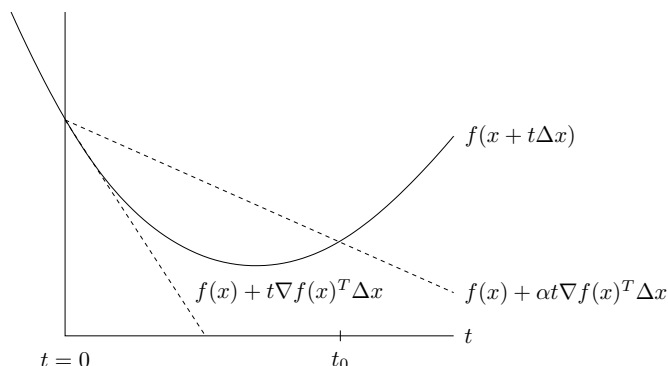


Figure 9.1 *Backtracking line search.* The curve shows f , restricted to the line over which we search. The lower dashed line shows the linear extrapolation of f , and the upper dashed line has a slope a factor of α smaller. The backtracking condition is that f lies below the upper dashed line, *i.e.*, $0 \leq t \leq t_0$.

The best values of β and γ vary by application; a good starting point is to take $\beta = 0.5$ and $\gamma = 0.05$.

Gradient descent with constraints

If there are constraints which we can specify through a restriction of the optimization variable to a set $\mathcal{C} \subset \mathbb{R}^N$,

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{x} \in \mathcal{C}, \quad (2)$$

then there is an analogous algorithm called **projected gradient descent**. For our discussion here, we will assume that the set \mathcal{C} is convex as is the functional f .

¹They are using ‘ α ’ in the figure label instead of the ‘ γ ’ we are using in these notes.

The projected gradient descent algorithm makes use of the “closest point” operator $P_{\mathcal{C}} : \mathbb{R}^N \rightarrow \mathbb{R}^N$,

$$P_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

As the name suggests, this takes an arbitrary vector in \mathbb{R}^N and returns the point in \mathcal{C} that is closest. If \mathcal{C} is closed and convex, then the closest point exists and is unique. For some sets \mathcal{C} it is easy to compute $P_{\mathcal{C}}(\mathbf{x})$, for others it is extremely hard.

Examples:

- \mathcal{C} is a K -dimensional subspace spanned by basis vectors that are the columns of the $N \times K$ matrix \mathbf{A} . Then

$$P_{\mathcal{C}}(\mathbf{x}) = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{x}.$$

If we have an orthobasis for the subspace concatenated as the columns of \mathbf{U} , this simplifies to

$$P_{\mathcal{C}}(\mathbf{x}) = \mathbf{U} \mathbf{U}^T \mathbf{x}.$$

For example, if $N = 4$, and

$$\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^4 : x[2] = 0, x[3] = 0\},$$

then

$$P_{\mathcal{C}}(\mathbf{x}) = \begin{bmatrix} x[1] \\ 0 \\ 0 \\ x[4] \end{bmatrix}.$$

- \mathcal{C} is the unit ball for the ℓ_2 norm:

$$\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\}.$$

Then

$$P_{\mathcal{C}}(\mathbf{x}) = \begin{cases} \mathbf{x}, & \|\mathbf{x}\|_2 \leq 1, \\ \frac{\mathbf{x}}{\|\mathbf{x}\|_2}, & \|\mathbf{x}\|_2 > 1. \end{cases}$$

- \mathcal{C} is the positive orthant:

$$\mathcal{C} = \{\mathbf{x} : x[n] \geq 0, n = 1, \dots, N\}.$$

Then

$$\hat{\mathbf{x}} = P_{\mathcal{C}}(\mathbf{x}) \quad \text{obeys} \quad \hat{x}[n] = \max(x[n], 0).$$

- \mathcal{C} is the set of $N \times N$ symmetric positive semi-definite matrices:

$$\mathcal{C} = \{\mathbf{X} : \mathbf{X} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T, \lambda_n \geq 0\}.$$

If \mathbf{X} is an arbitrary symmetric matrix with eigenvalue decomposition

$$\mathbf{X} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \sum_{n=1}^N \lambda_n \mathbf{v}_n \mathbf{v}_n^T,$$

then²

$$P_{\mathcal{C}}(\mathbf{X}) = \sum_{n=1}^N \max(\lambda_n, 0) \mathbf{v}_n \mathbf{v}_n^T.$$

If \mathbf{X} is an arbitrary (not necessarily symmetric) matrix, then we form

$$\mathbf{X}' = \frac{1}{2}(\mathbf{X} + \mathbf{X}^T)$$

with eigenvalue decomposition $\mathbf{X}' = \mathbf{V}'\mathbf{\Lambda}'\mathbf{V}'^T$ and take

$$P_{\mathcal{C}}(\mathbf{X}) = \sum_{n=1}^N \max(\lambda'_n, 0) \mathbf{v}'_n \mathbf{v}'_n{}^T.$$

²That this is indeed the closest point isn't at all obvious, but we will forgo the proof here.

If we do indeed have a tractable projection operator for \mathcal{C} , then we can solve (2) using the iterations:

$$\begin{aligned}\mathbf{z}_k &= \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \\ \mathbf{x}_{k+1} &= P_{\mathcal{C}}(\mathbf{z}_k).\end{aligned}$$

The convergence properties of this algorithm, with the appropriate choice of step size α_k , are very similar to those for unconstrained gradient descent.

Linear constraints. When \mathcal{C} is a subspace, then (as we saw above) the projection operator is itself linear. Thus in this case

$$\begin{aligned}\mathbf{x}_{k+1} &= P_{\mathcal{C}}(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)) \\ &= P_{\mathcal{C}}(\mathbf{x}_k) - \alpha_k P_{\mathcal{C}}(\nabla f(\mathbf{x}_k)) \\ &= \mathbf{x}_k - \alpha_k \mathbf{g}_k, \quad \mathbf{g}_k = P_{\mathcal{C}}(\nabla f(\mathbf{x}_k)),\end{aligned}$$

where the last equality follows from the fact that \mathbf{x}_k is, by construction, already in \mathcal{C} . For optimization with linear constraints, then, we can compute the gradient, project it, then choose a step size using exact or backtracking line search.

Interior point iterations

Another approach to handling constraints is to keep iterates in the “interior” of the set \mathcal{C} . We say \mathbf{x} is an **interior point** for \mathcal{C} if there exists an $\epsilon > 0$ such that

$$\mathbf{x} + \epsilon \mathbf{h} \in \mathcal{C}, \quad \text{for all } \mathbf{h} \in \mathbb{R}^D \text{ with } \|\mathbf{h}\|_2 \leq 1.$$

This basically means that the set \mathcal{C} has non-zero volume. Not all sets \mathcal{C} have an interior point — for example a subspace of dimension $K < N$ does not have any interior points.

Examples:

- If \mathcal{C} is the positive orthant, then all \mathbf{x} with $x[n] > 0$, $n = 1, \dots, N$ are interior points.
- If \mathcal{C} is the set of $N \times N$ symmetric positive semi-definite matrices, then all symmetric positive definite matrices (symmetric matrices that have eigenvalues that are strictly positive) are interior points.

If we start at an interior point \mathbf{x}_0 , we can solve (2) by using gradient descent with α_k chosen using backtracking line search with the additional constraint that

$$\mathbf{x}_k + t\mathbf{d}_k \in \text{Interior}(\mathcal{C}).$$

This type of algorithm can be very effective for small problems, but tends to converge very slowly for large problems when the solution is on the boundary of \mathcal{C} (where it tends to be, if the constraints are changing the problem in any way). Indeed, there is an entire class of interior point algorithms that are far more graceful, using smooth “barrier functions” to make the problem better posed.

A particular place where interior point iterations are interesting, and one that is relevant for Homework 10, is when

$$\mathcal{C} = \mathcal{X} \cap \mathcal{S},$$

where \mathcal{S} is a subspace (e.g. constraining certain optimization variables to be zero), and \mathcal{X} is a convex set with a non-empty interior³ (e.g. the set of symmetric positive semi-definite matrices). Now, from

³We really only need that \mathcal{X} has a non-empty interior relative to \mathcal{S} .

an initial point $\mathbf{x}_0 \in \mathcal{C}$, we can proceed via

$$\begin{aligned}\mathbf{g}_k &= \mathbf{P}_{\mathcal{S}}(\nabla f(\mathbf{x}_k)), \\ \mathbf{x}_{k+1} &= \mathbf{x}_k - \alpha_k \mathbf{g}_k,\end{aligned}$$

where α_k is chosen small enough so that $\mathbf{x}_{k+1} \in \text{Interior}(\mathcal{X})$.

References

- [BV04] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.