

Algorithms for constrained optimization

There are many, many constrained optimization algorithms, each tuned to the particulars of different classes of problems. We will look at the basics that underlie some of the more modern techniques. We will see that the concept of duality both helps us understand how these algorithms work, and gives us a way of determining when we are close to the solution.

We will describe several techniques. A recurring theme will be to replace the constrained program with an equivalent unconstrained program (or a series of unconstrained programs).

Eliminating equality constraints

The first approach is not so much an algorithm as a “trick” that lets us sometimes avoid even thinking about the constraints. Programs with *linear* equality constraints can always be written as programs without, and if there are no inequality constraints then the new program is unconstrained. To see this, suppose we are solving

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{b}.$$

Let \mathbf{x}_0 be any point satisfying $\mathbf{A}\mathbf{x}_0 = \mathbf{b}$. Then any feasible \mathbf{x} can be written as $\mathbf{x} = \mathbf{x}_0 + \mathbf{h}$, where $\mathbf{h} \in \text{Null}(\mathbf{A})$. Note that $\text{Null}(\mathbf{A})$ is a linear subspace of dimension $K = N - \text{rank}(\mathbf{A})$. If \mathbf{Q} is a basis for this space, we can write any $\mathbf{h} \in \text{Null}(\mathbf{A})$ as $\mathbf{h} = \mathbf{Q}\mathbf{w}$. Using this we can re-write the program above as

$$\underset{\mathbf{w} \in \mathbb{R}^K}{\text{minimize}} \quad f(\mathbf{x}_0 + \mathbf{Q}\mathbf{w}).$$

Sometimes this method can be very helpful, but note that computing \mathbf{x}_0 and \mathbf{Q} could potentially be expensive.

Projected gradient descent

Now suppose that we wish to solve the constrained optimization problem

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

where f is a differentiable convex function and \mathcal{C} is a convex set in \mathbb{R}^N . Another way to express this problem is as the unconstrained problem

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

where $\iota_{\mathcal{C}}$ denotes the indicator function for the set \mathcal{C} . We have previously encountered this idea in the context of duality, but in terms of suggesting practical algorithms, this has some obvious shortcomings. Namely, since $\iota_{\mathcal{C}}(\mathbf{x})$ is non-differentiable, we cannot apply gradient-based methods to solving (1).

However, we have encountered some algorithms for minimizing non-smooth convex functions. One that might seem particularly well-aligned with (1) is the proximal gradient method. Recall that proximal gradient method applies when our objective function consists of the sum of a smooth term (in this case, f) and a nonsmooth term (in this case, $\iota_{\mathcal{C}}$), resulting in the core iteration of

$$\mathbf{x}_{k+1} = \text{prox}_{\alpha_k \iota_{\mathcal{C}}}(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)).$$

This would yield a tractable algorithm with provable guarantees if we can compute $\text{prox}_{\alpha_k \iota_{\mathcal{C}}}$. So what does $\text{prox}_{\alpha_k \iota_{\mathcal{C}}}$ look like? Note that

$$\begin{aligned} \text{prox}_{\alpha_k \iota_{\mathcal{C}}}(\mathbf{z}) &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left(\iota_{\mathcal{C}}(\mathbf{x}) + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{z}\|_2^2 \right) \\ &= \arg \min_{\mathbf{x} \in \mathcal{C}} \|\mathbf{x} - \mathbf{z}\|_2^2 \\ &= \mathcal{P}_{\mathcal{C}}(\mathbf{z}), \end{aligned}$$

where $\mathcal{P}_{\mathcal{C}}(\mathbf{z})$ denotes the **projection** of \mathbf{z} onto the set \mathcal{C} . Note that this holds for any $\alpha_k > 0$.

Thus, the core iteration of the proximal gradient method is equivalent to

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

That is, at each iteration we take a gradient step on f and then re-project onto the constraint set. Notice that for any $k \geq 1$, we are always guaranteed that \mathbf{x}_k is feasible.

This algorithm is usually called **projected gradient descent**. It is a very simple (but often effective) method for solving constrained optimization problems when the projection onto the constraint set \mathcal{C} can be computed efficiently. Note, however, that this is not always the case – sometimes computing this projection itself requires solving a challenging optimization problem.

Of course projected gradient descent inherits the convergence guarantees for the proximal gradient method. In particular, if f is L -smooth, then we know that the iterates obey

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

Other convergence results mirror those for unconstrained gradient descent. In particular, if f is smooth and strongly convex, then projected gradient descent has linear convergence. A full set of results (along with detailed proofs) can be found in [Lan20, Chap. 3].

Example: Least-squares with positivity constraints

Suppose we want to solve

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

This is a case where the projection onto the constraint set is relatively simple. It is easy to argue that the projection onto the set of all positive vectors is to simply just set all of the negative entries to zero. The projected gradient descent iteration is then

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

where

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

Barrier methods

Another popular and even more flexible approach are **barrier methods**. These can be thought of as again replacing our constrained problem with the unconstrained one in (1), but rather than attempting to minimize (1) directly, we instead solve a slight perturbation of this problem. In particular, we replace the indicator function $\iota_{\mathcal{C}}$ with some function b such that $\text{dom } b = \mathcal{C}$ and $b(\mathbf{x}) \rightarrow \infty$ as \mathbf{x} approaches the boundary of \mathcal{C} .

To make this concrete, consider the constrained program

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) \quad \text{subject to} \quad g_m(\mathbf{x}) \leq 0, \quad m = 1, \dots, M.$$

In a barrier method we replace this with the unconstrained program

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) + \sum_{m=1}^M B(g_m(\mathbf{x})),$$

where $\text{dom } B = \mathbb{R}_-$ and $B(x) \rightarrow \infty$ as $x \rightarrow 0$ from the left. Again, unless B is the indicator function, the new program is an approximation to the original.

An interesting choice is $B(x) = -\frac{1}{\tau} \log(-x)$. This particular barrier function has the properties:

- You can analyze the number of Newton iterations needed for convergence for many f of interest (using self-concordance). This is done very nicely in Chapters 9 and 11 of [BV04].
- The solution¹ $\mathbf{x}^*(\tau)$ of

$$\underset{\mathbf{x}}{\text{minimize}} \quad \tau f(\mathbf{x}) - \sum_{m=1}^M \log(-g_m(\mathbf{x}))$$

can be used to generate a dual-feasible point (and hence a primal-dual gap certificate), and can be related to the KKT conditions for the original program.

To appreciate the second point above, we start by taking the gradient of the objective function above and setting it equal to zero. We see that

$$\tau \nabla f(\mathbf{x}^*(\tau)) + \sum_{m=1}^M -\frac{1}{g_m(\mathbf{x}^*(\tau))} \nabla g_m(\mathbf{x}^*(\tau)) = \mathbf{0}.$$

So if we take

$$\lambda_m^*(\tau) = -\frac{1}{\tau g_m(\mathbf{x}^*(\tau))}, \quad m = 1, \dots, M,$$

we have $\lambda_m^*(\tau) \geq 0$ and

$$f(\mathbf{x}^*(\tau)) + \sum_{m=1}^M \lambda_m^*(\tau) \nabla g_m(\mathbf{x}^*(\tau)) = \mathbf{0}.$$

¹We have multiplied the objective by τ to make some of what follows a little easier to express.

Since $\mathbf{x}^*(\tau)$ is primal feasible, the only KKT condition we are missing is complementary slackness – we have replaced the condition

$$\lambda_m^* g_m(\mathbf{x}^*) = 0, \quad \text{with} \quad \lambda_m^*(\tau) g_m(\mathbf{x}^*(\tau)) = -1/\tau.$$

So if we set τ to be increasingly large, we obtain points that satisfy an increasingly tight approximation to the KKT conditions.

With this choice of $\boldsymbol{\lambda}^*(\tau)$, we can also easily compute the dual of the original program:

$$\begin{aligned} d(\boldsymbol{\lambda}^*(\tau)) &= \inf_{\mathbf{x}} \left(f(\mathbf{x}) + \sum_{m=1}^M \lambda_m^*(\tau) g_m(\mathbf{x}) \right) \\ &= f(\mathbf{x}^*(\tau)) + \sum_{m=1}^M \lambda_m^*(\tau) g_m(\mathbf{x}^*(\tau)) \\ &= f(\mathbf{x}^*(\tau)) - M/\tau. \end{aligned}$$

Hence, we know that if \mathbf{x}^* is a solution to the original program, then

$$f(\mathbf{x}^*(\tau)) - f(\mathbf{x}^*) \leq f(\mathbf{x}^*(\tau)) - d(\boldsymbol{\lambda}^*(\tau)) \leq M/\tau.$$

Thus, we know that solving the log-barrier problem gets us within M/τ of the optimal of the original primal objective.

A full discussion of log barrier methods, including some fundamental complexity analysis, can be found in [BV04, Chap. 11]. One interesting theoretical result there is that, with a reasonable way of adjusting τ (multiplying it by 10 at every iteration, for example), the number of log-barrier iterations to make the value of the barrier functional $f(\mathbf{x}^*(\tau))$ agree with the minimal value p^* to the original constrained problem to some precision. The upshot is that there is a very close match after $\sim \sqrt{M}$ iterations. This means that in theory, solving a

constrained problem is roughly as expensive as solving \sqrt{M} unconstrained problems. In practice, it is actually much cheaper – standard log barrier iterations take maybe 20–50 iterations to produce good results.

Primal dual interior point methods

These are closely related to log barrier algorithms, but they take a more direct approach towards “solving” the KKT conditions. The general idea is to treat the KKT conditions like a set of nonlinear equations, and solve them using Newton’s method.

We start with the same set of relaxed KKT conditions² we used with log barrier:

$$\mathbf{r}_\tau(\mathbf{x}, \boldsymbol{\lambda}) = \begin{bmatrix} \nabla f(\mathbf{x}) + \sum_m \lambda_m \nabla g_m(\mathbf{x}) \\ -\lambda_1 g_1(\mathbf{x}) - 1/\tau \\ \vdots \\ -\lambda_m g_m(\mathbf{x}) - 1/\tau \end{bmatrix}$$

If we find \mathbf{x} and $\boldsymbol{\lambda}$ such that the $N + M$ -vector $\mathbf{r}_\tau(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{0}$, then we know we have found the same $\mathbf{x}^*(\tau)$, $\boldsymbol{\lambda}^*(\tau)$ that solve the log barrier problem.

Primal-dual interior point methods take Newton steps to try to make $\mathbf{r}_\tau = \mathbf{0}$, but they adjust τ at every step. The Newton step is characterized by

$$\mathbf{r}_\tau(\mathbf{x} + \delta\mathbf{x}, \boldsymbol{\lambda} + \delta\boldsymbol{\lambda}) \approx \mathbf{r}_\tau(\mathbf{x}, \boldsymbol{\lambda}) + \mathbf{J}_r(\mathbf{x}, \boldsymbol{\lambda}) \begin{bmatrix} \delta\mathbf{x} \\ \delta\boldsymbol{\lambda} \end{bmatrix} = 0,$$

²We are again only considering inequality constraints; it is straightforward to modify everything we say here to include linear equality constraints.

where $\mathbf{J}_{\mathbf{r}_\tau}(\mathbf{x}, \boldsymbol{\lambda})$ is the **Jacobian** matrix for the vector-valued function \mathbf{r}_τ given by

$$\begin{bmatrix} \nabla^2 f(\mathbf{x}) + \sum_m \lambda_m \nabla^2 g_m(\mathbf{x}) & \nabla g_1(\mathbf{x}) & \nabla g_2(\mathbf{x}) & \cdots & \nabla g_M(\mathbf{x}) \\ -\lambda_1 \nabla g_1(\mathbf{x})^\top & -g_1(\mathbf{x}) & 0 & \cdots & 0 \\ -\lambda_2 \nabla g_2(\mathbf{x})^\top & 0 & -g_2(\mathbf{x}) & \cdots & 0 \\ \vdots & & & \ddots & \\ -\lambda_M \nabla g_M(\mathbf{x})^\top & 0 & 0 & \cdots & -g_M(\mathbf{x}) \end{bmatrix}.$$

The update direction is

$$\begin{bmatrix} \delta \mathbf{x} \\ \delta \boldsymbol{\lambda} \end{bmatrix} = -\mathbf{J}_{\mathbf{r}_\tau}^{-1} \mathbf{r}_\tau(\mathbf{x}, \boldsymbol{\lambda}).$$

With this direction, we can perform a line search. The parameter τ is updated at every step (getting larger) based on an estimate of the duality gap.

A key feature of this type of primal dual method is that the iterates \mathbf{x}_k and $\boldsymbol{\lambda}_k$ do not have to be feasible (although they of course become feasible in the limit).

Details on this particular algorithm, along with a full convergence analysis, can be found in [BV04, Chap. 11.7].

References

- [BV04] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [Lan20] G. Lan. *First-order and Stochastic Optimization Methods for Machine Learning*. Springer, 2020.