

Alternating direction primal dual methods

We will now focus on a class of algorithms that work by fixing the dual variables and updating the primal variables \mathbf{x} , then fixing the primals and updating the dual variables $\boldsymbol{\lambda}, \boldsymbol{\nu}$. An excellent source for this material is [BPC⁺10]. In fact, what follows here is basically a summary of the first 12 pages of that paper.

We have seen that when we have strong duality (which we will assume throughout), the optimal value of the primal program is equal to the optimal value of the dual program. That is, if $\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*$ are primal/dual optimal points,

$$\begin{aligned} f(\mathbf{x}^*) &= d(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) \\ &= \inf_{\mathbf{x} \in \mathbb{R}^N} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*), \end{aligned}$$

where \mathcal{L} is the Lagrangian

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f(\mathbf{x}) + \sum_{m=1}^M \lambda_m g_m(\mathbf{x}) + \langle \mathbf{A}\mathbf{x} - \mathbf{b}, \boldsymbol{\nu} \rangle.$$

If $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$ has only one minimizer,¹ then we can recover the primal optimal solution \mathbf{x}^* from the dual-optimal solution $\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*$ by solving the *unconstrained* program

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*).$$

“Alternating” methods search for a saddle point of the Lagrangian by fixing the dual variables $\boldsymbol{\lambda}_k, \boldsymbol{\nu}_k$, minimizing $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k, \boldsymbol{\nu}_k)$ with respect to \mathbf{x} , then updating the Lagrange multipliers.

¹Which is the case when f is strictly convex, and in many other situations.

To start, we will base our discussion on **equality constrained** problems. Incorporating inequality constraints will be natural after we have developed things a bit.

Dual ascent

We want to solve

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

We will assume that the domain of f is all of \mathbb{R}^N ; again, things are easily modified if this is any open set. The Lagrangian is

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\nu}) = f(\mathbf{x}) + \langle \mathbf{A}\mathbf{x} - \mathbf{b}, \boldsymbol{\nu} \rangle,$$

and the dual function is

$$\begin{aligned} d(\boldsymbol{\nu}) &= \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}) \\ &= \inf_{\mathbf{x}} f(\mathbf{x}) + \langle \mathbf{A}\mathbf{x}, \boldsymbol{\nu} \rangle - \langle \mathbf{b}, \boldsymbol{\nu} \rangle \\ &= -\sup_{\mathbf{x}} \left(\langle \mathbf{x}, -\mathbf{A}^T \boldsymbol{\nu} \rangle - f(\mathbf{x}) \right) - \langle \mathbf{b}, \boldsymbol{\nu} \rangle \\ &= -f^*(-\mathbf{A}^T \boldsymbol{\nu}) - \langle \mathbf{b}, \boldsymbol{\nu} \rangle \end{aligned}$$

and the dual problem is

$$\underset{\boldsymbol{\nu} \in \mathbb{R}^N}{\text{maximize}} \quad d(\boldsymbol{\nu}).$$

Consider for a moment the problem of maximizing the dual. A reasonable thing to do would be some kind of gradient ascent:²

$$\boldsymbol{\nu}_{k+1} = \boldsymbol{\nu}_k + \alpha_k \nabla d(\boldsymbol{\nu}_k),$$

²“Ascent” instead of “descent” because g is concave instead of convex.

where α_k is some appropriate step size. The gradient of d (with respect to $\boldsymbol{\nu}$) at a point $\boldsymbol{\nu}_0$ is

$$\nabla d(\boldsymbol{\nu}_0) = \nabla \inf_{\mathbf{x}} (f(\mathbf{x}) + \langle \boldsymbol{\nu}_0, \mathbf{A}\mathbf{x} - \mathbf{b} \rangle),$$

and so if $\mathbf{x}^+ = \arg \min_{\mathbf{x}} (f(\mathbf{x}) + \langle \boldsymbol{\nu}_0, \mathbf{A}\mathbf{x} - \mathbf{b} \rangle)$, then

$$\begin{aligned} \nabla d(\boldsymbol{\nu}_0) &= \nabla_{\boldsymbol{\nu}} (f(\mathbf{x}^+) + \langle \boldsymbol{\nu}_0, \mathbf{A}\mathbf{x}^+ - \mathbf{b} \rangle) \\ &= \mathbf{A}\mathbf{x}^+ - \mathbf{b}. \end{aligned}$$

This leads naturally to:

The **dual ascent** algorithm consists of the iteration

$$\begin{aligned} \mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}_k) \\ \boldsymbol{\nu}_{k+1} &= \boldsymbol{\nu}_k + \alpha_k (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{b}) \end{aligned}$$

that is repeated until some convergence criteria is met.

This algorithm “works” under certain assumptions on f (that translate to different assumptions on the dual d). In particular, we need $\mathcal{L}(\mathbf{x}, \boldsymbol{\nu})$ to be bounded for every $\boldsymbol{\nu}$, otherwise the primal update $\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}_k)$ can fail.

That the Lagrangian is bounded below for every choice of $\boldsymbol{\nu}$ is far from a given. For example, a program of the form

$$\underset{\mathbf{x}}{\text{minimize}} \quad \langle \mathbf{x}, \mathbf{c} \rangle \quad \text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{b}$$

will have Lagrangian

$$\begin{aligned} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}) &= \langle \mathbf{x}, \mathbf{c} \rangle + \boldsymbol{\nu}^T (\mathbf{A}\mathbf{x} - \mathbf{b}) \\ &= (\mathbf{c} + \mathbf{A}^T \boldsymbol{\nu})^T \mathbf{x} - \boldsymbol{\nu}^T \mathbf{b}, \end{aligned}$$

which is unbounded below since the first term is linear in \mathbf{x} .

Of course, this algorithm is nicest when we can solve the unconstrained primal update problem efficiently.

Dual decomposition

Dual ascent is simple, and it is pretty much as old an idea as convex optimization itself. But it has a key feature that makes it very attractive (at least as a starting point) for modern computing platforms.

If the objective functional f is **separable**, we can also separate (i.e., parallelize) the primal update. Suppose we can write $f(\mathbf{x})$ as a sum

$$f(\mathbf{x}) = \sum_{i=1}^B f_i(\mathbf{x}^{(i)}),$$

where the $\mathbf{x}^{(i)} \in \mathbb{R}^{N_i}$ ($\sum_i N_i = N$) are a partition of \mathbf{x} . We can partition the columns of \mathbf{A} in the same way,

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}^{(1)} & \mathbf{A}^{(2)} & \dots & \mathbf{A}^{(B)} \end{bmatrix}$$

so that

$$\mathbf{A}\mathbf{x} = \sum_{i=1}^B \mathbf{A}^{(i)} \mathbf{x}^{(i)}.$$

Notice that even with this assumption, the primal optimization program itself is not separable, as the constraints $\mathbf{A}\mathbf{x} = \mathbf{b}$ tie all of the blocks together – there is still just a single right-hand side vector \mathbf{b} .

The Lagrangian in this case is also separable:

$$\begin{aligned}
\mathcal{L}(\mathbf{x}, \boldsymbol{\nu}) &= \sum_{i=1}^B f_i(\mathbf{x}^{(i)}) + \boldsymbol{\nu}^T \left(\sum_{i=1}^B \mathbf{A}^{(i)} \mathbf{x}^{(i)} \right) - \boldsymbol{\nu}^T \mathbf{b} \\
&= \sum_{i=1}^B \left[f_i(\mathbf{x}^{(i)}) + \boldsymbol{\nu}^T \mathbf{A}^{(i)} \mathbf{x}^{(i)} - \frac{1}{B} \boldsymbol{\nu}^T \mathbf{b} \right] \\
&= \sum_{i=1}^B \mathcal{L}_i(\mathbf{x}^{(i)}, \boldsymbol{\nu}).
\end{aligned}$$

For fixed $\boldsymbol{\nu}_k$, the primal update

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\nu}_k)$$

can be separated into B independent updates:

$$\mathbf{x}_{k+1}^{(i)} = \arg \min_{\mathbf{x}^{(i)}} \mathcal{L}_i(\mathbf{x}^{(i)}, \boldsymbol{\nu}_k)$$

You can imagine broadcasting the current dual iterate to B different processors; they each compute their piece of the primal update $\mathbf{x}_{k+1}^{(i)}$, and then the results are gathered centrally to compute the dual update. Notice also, though, that the hard part of the dual update, computing $\mathbf{A}\mathbf{x}_{k+1}$, can also be done in a decentralized manner.

The Method of Multipliers and Augmented Lagrangians

The method of multipliers (MoM) is the same idea as dual ascent, but we smooth out (augment) the Lagrangian to make the primal update more robust.

It should be clear that

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

and

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 \quad \text{subject to} \quad \mathbf{Ax} = \mathbf{b}$$

have exactly the same set of solutions for all $\rho \geq 0$.

The Lagrangian for the second program is

$$\mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\nu}) = f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \boldsymbol{\nu}^T(\mathbf{Ax} - \mathbf{b}).$$

This is called the **augmented Lagrangian** of the original problem.

Adding the quadratic term is nice – it makes (under mild conditions on f with respect to \mathbf{A}) the primal update minimization well-posed (i.e., makes the dual differentiable).

Notice that the Lagrange multipliers $\boldsymbol{\nu}$ appear in exactly the same way in the augmented Lagrangian as they do in the regular Lagrangian, so the dual update does not change.

The resulting algorithm is called the **method of multipliers**; we iterate

$$\begin{aligned}\mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \mathcal{L}_{\rho}(\mathbf{x}, \boldsymbol{\nu}_k) \\ \boldsymbol{\nu}_{k+1} &= \boldsymbol{\nu}_k + \rho(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{b})\end{aligned}$$

until some convergence criteria is met.

As a bonus, we now have a principled way of selecting the step size for the dual update – just use ρ . To see why this makes sense, recall the KKT conditions for \mathbf{x}^* and $\boldsymbol{\nu}^*$ to be a solution:

$$\mathbf{A}\mathbf{x}^* = \mathbf{b}, \quad \nabla f(\mathbf{x}^*) + \mathbf{A}^T \boldsymbol{\nu}^* = \mathbf{0}.$$

With ρ as the step size, we have

$$\begin{aligned}\mathbf{0} &= \nabla \mathcal{L}_{\rho}(\mathbf{x}_{k+1}, \boldsymbol{\nu}_k), \quad (\text{since } \mathbf{x}_{k+1} \text{ is a minimizer}), \\ &= \nabla f(\mathbf{x}_{k+1}) + \mathbf{A}^T (\boldsymbol{\nu}_k + \rho(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{b})) \\ &= \nabla f(\mathbf{x}_{k+1}) + \mathbf{A}^T \boldsymbol{\nu}_{k+1}.\end{aligned}$$

So the dual update maintains the second optimality condition at every step.

The MoM has much better convergence properties than dual ascent, but it is no longer separable. The algorithm we look at next, the alternating direction method of multipliers (ADMM), will build on this idea in such a way that we do have a type of dual decomposition for the smoothed Lagrangian. In addition, it can be easily modified to incorporate certain kinds of inequality constraints.

References

- [BPC⁺10] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2010.