

Alternating direction method of multipliers (ADMM)

ADMM extends the method of multipliers in such a way that we get back some of the decomposability (i.e., ability to parallelize) of standard dual ascent algorithms. It also gives us a flexible framework for incorporating many types of convex constraints, though we will again focus on linear equality constraints to start. See [BPC⁺10] for a more complete discussion.

ADMM **splits** the optimization variables into two parts, \mathbf{x} and \mathbf{z} , and solves programs of the form

$$\underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{z}) \quad \text{subject to} \quad \mathbf{Ax} + \mathbf{Bz} = \mathbf{c}.$$

The basic idea is to rotate through 3 steps:

1. Minimize the (augmented) Lagrangian over \mathbf{x} with \mathbf{z} and the Lagrange multipliers $\boldsymbol{\nu}$ fixed.
2. Minimize the (augmented) Lagrangian over \mathbf{z} with \mathbf{x} and $\boldsymbol{\nu}$ fixed.
3. Update the Lagrange multipliers using gradient ascent as before.

If the splitting is done in a careful manner, it can happen that each of the subproblems above can be easily computed. We can also handle general convex constraints (more on this later).

To make the three steps above more explicit: the augmented Lagrangian is

$$\mathcal{L}_\rho(\mathbf{x}, \mathbf{z}, \boldsymbol{\nu}) = f(\mathbf{x}) + g(\mathbf{z}) + \boldsymbol{\nu}^\top (\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}) + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}\|_2^2,$$

and the general ADMM iteration is

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

The only real difference between ADMM and MoM is the we are splitting the primal minimization into two parts instead of optimizing over (\mathbf{x}, \mathbf{z}) jointly.

Scaled form.

We can write the ADMM iterations in a more convenient form by substituting

$$\boldsymbol{\mu} = \frac{1}{\rho} \boldsymbol{\nu}.$$

By “completing the square” we have that

$$\boldsymbol{\nu}^T(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c}\|_2^2 = \frac{\rho}{2} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c} + \boldsymbol{\mu}\|_2^2 - \frac{\rho}{2} \|\boldsymbol{\mu}\|_2^2,$$

and so we can write:

ADMM:

$$\begin{aligned}\mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \left(f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z}_k - \mathbf{c} + \boldsymbol{\mu}_k\|_2^2 \right) \\ \mathbf{z}_{k+1} &= \arg \min_{\mathbf{z}} \left(g(\mathbf{z}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}\mathbf{z} - \mathbf{c} + \boldsymbol{\mu}_k\|_2^2 \right) \\ \boldsymbol{\mu}_{k+1} &= \boldsymbol{\mu}_k + \mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}\mathbf{z}_{k+1} - \mathbf{c}\end{aligned}$$

Example: The LASSO

Recall the LASSO:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \tau \|\mathbf{x}\|_1.$$

Taking

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 \quad \text{and} \quad g(\mathbf{z}) = \tau \|\mathbf{z}\|_1,$$

we can rewrite this in ADMM form as

$$\underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{z}) \quad \text{subject to} \quad \mathbf{x} - \mathbf{z} = \mathbf{0}.$$

The \mathbf{x} update is

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \left(\frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \frac{\rho}{2} \|\mathbf{x} - \mathbf{z}_k + \boldsymbol{\mu}_k\|_2^2 \right).$$

With both \mathbf{z}_k and $\boldsymbol{\mu}_k$ fixed, this is a regularized least-squares problem and is equivalent to:

$$\min_{\mathbf{x}} \left\| \begin{bmatrix} \mathbf{A} \\ \sqrt{\rho} \mathbf{I} \end{bmatrix} \mathbf{x} - \begin{bmatrix} \mathbf{b} \\ \sqrt{\rho}(\mathbf{z}_k - \boldsymbol{\mu}_k) \end{bmatrix} \right\|_2^2.$$

This problem has a closed-form solution:

$$\begin{aligned} \mathbf{x}_{k+1} &= \left(\mathbf{A}^T \mathbf{A} + \rho \mathbf{I} \right)^{-1} \left[\mathbf{A}^T \quad \sqrt{\rho} \mathbf{I} \right] \begin{bmatrix} \mathbf{b} \\ \sqrt{\rho}(\mathbf{z}_k - \boldsymbol{\mu}_k) \end{bmatrix} \\ &= \left(\mathbf{A}^T \mathbf{A} + \rho \mathbf{I} \right)^{-1} \left(\mathbf{A}^T \mathbf{b} + \rho(\mathbf{z}_k - \boldsymbol{\mu}_k) \right) \end{aligned}$$

The \mathbf{z} update problem is:

$$\underset{\mathbf{z}}{\text{minimize}} \quad \tau \|\mathbf{z}\|_1 + \frac{\rho}{2} \|\mathbf{z} - \mathbf{x}_{k+1} - \boldsymbol{\mu}_k\|_2^2.$$

You may recognize this: it is the proximal operator for the ℓ_1 -norm, which as we have seen before has closed form solution:

$$\mathbf{z}_{k+1} = T_{\tau/\rho}(\mathbf{x}_{k+1} + \boldsymbol{\mu}_k),$$

where $T_\lambda(\cdot)$ is the term-by-term soft-thresholding operator,

$$(T_\lambda(\mathbf{v}))[n] = \begin{cases} v[n] - \lambda, & v[n] > \lambda, \\ 0, & |v[n]| \leq \lambda, \\ v[n] + \lambda, & v[n] < -\lambda. \end{cases}$$

To summarize:

ADMM iterations for the LASSO

$$\begin{aligned} \mathbf{x}_{k+1} &= \left(\mathbf{A}^T \mathbf{A} + \rho \mathbf{I}\right)^{-1} \left(\mathbf{A}^T \mathbf{b} + \rho(\mathbf{z}_k - \boldsymbol{\mu}_k)\right), \\ \mathbf{z}_{k+1} &= T_{\tau/\rho}(\mathbf{x}_{k+1} + \boldsymbol{\mu}_k), \\ \boldsymbol{\mu}_{k+1} &= \boldsymbol{\mu}_k + \mathbf{x}_{k+1} - \mathbf{z}_{k+1}. \end{aligned}$$

Convergence properties

We will state one convergence result. If the following two conditions hold:

1. f and g are closed, proper, and convex (i.e., their epigraphs are nonempty closed convex sets),
2. strong duality holds,

then

- $\mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{z}_k - \mathbf{c} \rightarrow \mathbf{0}$ as $k \rightarrow \infty$. That is, the primal iterates are asymptotically feasible.
- $f(\mathbf{x}_k) + g(\mathbf{z}_k) \rightarrow f(\mathbf{x}^*) + g(\mathbf{z}^*)$ as $k \rightarrow \infty$. That is, the value of the objective function approaches the optimal value asymptotically.
- $\boldsymbol{\nu}_k \rightarrow \boldsymbol{\nu}^*$ as $k \rightarrow \infty$, where $\boldsymbol{\nu}^*$ is a dual optimal point.

Under additional assumptions, we can also have convergence to a primal optimal point, $(\mathbf{x}_k, \mathbf{z}_k) \rightarrow (\mathbf{x}^*, \mathbf{z}^*)$ as $k \rightarrow \infty$.

See [BPC⁺10, Section 3.2] for further discussion and references.

Convex constraints

Using a technique that we have seen before, we can write the general program

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

where \mathcal{C} is a closed convex set, in ADMM form as

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{z}) \quad \text{subject to} \quad \mathbf{x} - \mathbf{z} = \mathbf{0},$$

where $g(\mathbf{z})$ is the indicator function for \mathcal{C} :

$$g(\mathbf{z}) = \begin{cases} 0, & \mathbf{z} \in \mathcal{C}, \\ \infty, & \mathbf{z} \notin \mathcal{C}. \end{cases}$$

Note that in this case, the \mathbf{z} update is a closest-point-to-a-convex-set problem. For fixed $\mathbf{v} \in \mathbb{R}^N$:

$$\arg \min_{\mathbf{z}} g(\mathbf{z}) + \frac{\rho}{2} \|\mathbf{z} - \mathbf{v}\|_2^2 = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{v}\|_2 = \mathcal{P}_{\mathcal{C}}(\mathbf{v}).$$

ADMM iteration for general convex constraints:

$$\begin{aligned} \mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \left(f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{x} - \mathbf{z}_k + \boldsymbol{\mu}_k\|_2^2 \right), \\ \mathbf{z}_{k+1} &= \mathcal{P}_{\mathcal{C}}(\mathbf{x}_{k+1} + \boldsymbol{\mu}_k), \\ \boldsymbol{\mu}_{k+1} &= \boldsymbol{\mu}_k + \mathbf{x}_{k+1} - \mathbf{z}_{k+1}. \end{aligned}$$

Of course, this algorithm is most attractive when we have a fast method for computing $\mathcal{P}_{\mathcal{C}}(\cdot)$.

Example: Basis Pursuit

As we have seen before, a good proxy for finding the sparsest solution to an underdetermined system of equations $\mathbf{Ax} = \mathbf{b}$ is to solve

$$\underset{\mathbf{x}}{\text{minimize}} \|\mathbf{x}\|_1 \quad \text{subject to} \quad \mathbf{Ax} = \mathbf{b}.$$

To put this in ADMM form, we are solving

$$\underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} f(\mathbf{x}) + g(\mathbf{z}) \quad \text{subject to} \quad \mathbf{x} - \mathbf{z} = \mathbf{0},$$

with

$$f(\mathbf{x}) = \|\mathbf{x}\|_1, \quad \text{and} \quad g(\mathbf{z}) = \begin{cases} 0, & \mathbf{Az} = \mathbf{b}, \\ \infty, & \text{otherwise.} \end{cases}$$

The projection onto $\mathcal{C} = \{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$ can be given in closed form using the pseudo-inverse c of \mathbf{A} as

$$\begin{aligned} \mathcal{P}_{\mathcal{C}}(\mathbf{v}) &= \mathbf{A}^\dagger(\mathbf{b} - \mathbf{Av}) + \mathbf{v} \\ &= (\mathbf{I} - \mathbf{A}^\top(\mathbf{AA}^\top)^{-1}\mathbf{A})\mathbf{v} + \mathbf{A}^\top(\mathbf{AA}^\top)^{-1}\mathbf{b}, \end{aligned}$$

where the last equality comes from $\mathbf{A}^\dagger = \mathbf{A}^\top(\mathbf{AA}^\top)^{-1}$ when \mathbf{A} has full row rank.

The updates in this case are

$$\begin{aligned} \mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \left(\|\mathbf{x}\|_1 + \frac{\rho}{2} \|\mathbf{x} - \mathbf{z}_k + \boldsymbol{\mu}_k\|_2^2 \right) \\ &= T_{1/\rho}(\mathbf{z}_k - \boldsymbol{\mu}_k) \\ \mathbf{z}_{k+1} &= (\mathbf{I} - \mathbf{A}^\top(\mathbf{AA}^\top)^{-1}\mathbf{A})(\mathbf{x}_{k+1} + \boldsymbol{\mu}_k) + \mathbf{A}^\top(\mathbf{AA}^\top)^{-1}\mathbf{b} \\ \boldsymbol{\mu}_{k+1} &= \boldsymbol{\mu}_k + \mathbf{x}_{k+1} - \mathbf{z}_{k+1}. \end{aligned}$$

Example: Linear programming

Consider the general linear program

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{c}^T \mathbf{x} \quad \text{subject to} \quad \mathbf{A} \mathbf{x} = \mathbf{b}, \quad \mathbf{x} \geq \mathbf{0},$$

where \mathbf{A} is an $M \times N$ matrix with full row rank.¹ We can put this in ADMM form by first eliminating the equality constraints, then introducing the indicator function for the non-negativity constraint.

Let \mathbf{Q} be an $N \times (N - M)$ matrix whose columns span $\text{Null}(\mathbf{A})$, and let \mathbf{x}_0 be any point such that $\mathbf{A} \mathbf{x}_0 = \mathbf{b}$. Then we can re-write the LP as

$$\underset{\mathbf{w}}{\text{minimize}} \quad \mathbf{c}^T (\mathbf{x}_0 + \mathbf{Q} \mathbf{w}) \quad \text{subject to} \quad \mathbf{x}_0 + \mathbf{Q} \mathbf{w} \geq \mathbf{0},$$

which we can write in ADMM form as

$$\underset{\mathbf{w}}{\text{minimize}} \quad \mathbf{c}^T \mathbf{x}_0 + \mathbf{c}^T \mathbf{Q} \mathbf{w} + g(\mathbf{z}) \quad \text{subject to} \quad \mathbf{Q} \mathbf{w} - \mathbf{z} = -\mathbf{x}_0,$$

where

$$g(\mathbf{z}) = \begin{cases} 0, & \mathbf{z} \geq \mathbf{0}, \\ \infty, & \text{otherwise.} \end{cases}$$

(We can drop the $\mathbf{c}^T \mathbf{x}_0$ from the objective since it does not depend on either of the optimization variables.)

Notice that when \mathbf{Q} has full column rank, the program

$$\underset{\mathbf{w}}{\text{minimize}} \quad \mathbf{v}^T \mathbf{w} + \frac{1}{2} \|\mathbf{Q} \mathbf{w} - \mathbf{y}\|_2^2,$$

¹The full row rank assumption is not at all essential; I am just making it to keep things streamlined.

has the closed-form solution

$$\mathbf{w}^* = (\mathbf{Q}^T \mathbf{Q})^{-1} (\mathbf{Q}^T \mathbf{y} - \mathbf{v}).$$

Also, the projection onto the non-negative orthant $\mathcal{C} = \{\mathbf{x} : \mathbf{x} \geq \mathbf{0}\}$ is

$$\mathcal{P}_{\mathcal{C}}(\mathbf{v}) = (\mathbf{v})_+, \quad \text{or} \quad (\mathcal{P}_{\mathcal{C}}(\mathbf{v}))_n = \begin{cases} v[n], & v[n] \geq 0, \\ 0, & v[n] < 0. \end{cases}$$

For the general linear program, then, the ADMM iterations are

$$\begin{aligned} \mathbf{w}_{k+1} &= \arg \min_{\mathbf{w}} \left(\frac{1}{\rho} \mathbf{c}^T \mathbf{Q} \mathbf{w} + \frac{1}{2} \|\mathbf{Q} \mathbf{w} - \mathbf{z}_k + \mathbf{x}_0 + \boldsymbol{\mu}_k\|_2^2 \right) \\ &= (\mathbf{Q}^T \mathbf{Q})^{-1} \left[\mathbf{Q}^T (\mathbf{z}_k - \mathbf{x}_0 - \boldsymbol{\mu}_k) - \frac{1}{\rho} \mathbf{Q}^T \mathbf{c} \right], \\ \mathbf{z}_{k+1} &= \mathcal{P}_{\mathcal{C}}(\mathbf{Q} \mathbf{w}_{k+1} + \mathbf{x}_0 + \boldsymbol{\mu}_k) \\ &= (\mathbf{Q} \mathbf{w}_{k+1} + \mathbf{x}_0 + \boldsymbol{\mu}_k)_+ \\ \boldsymbol{\mu}_{k+1} &= \boldsymbol{\mu}_k + \mathbf{Q} \mathbf{w}_{k+1} - \mathbf{z}_{k+1} + \mathbf{x}_0. \end{aligned}$$

Notice that especially when the columns of \mathbf{Q} are orthogonal, $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$, all of these steps are very simple.

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.