

Lecture 11

Distributed optimization and ADMM

Instructor: Prof. Gabriele Farina (✉ gfarina@mit.edu)*

In this last lecture on first-order methods, we will touch on an important aspect of optimization in machine learning: *distributed* optimization.

1 Setting

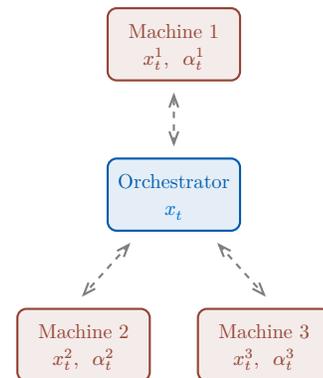
Consider again an empirical risk minimization problem, where we aim to minimize the average loss over a large training set. When the training set is too large to fit on a single machine, or privacy concerns require different parts of the data to be stored on different machines, it is common to distribute the optimization problem across multiple machines.

In this setting, we assume that we have m machines, each with access to a local dataset and a local function $f^j : \mathbb{R}^n \rightarrow \mathbb{R}$. The goal is to minimize the average of the local functions, that is,

$$f(x) = \frac{1}{m} \sum_{j=1}^m f^j(x),$$

where the functions f^j are differentiable.

In our distributed optimization model, each machine can communicate with a central **orchestrator**, which can send messages to all machines and collect messages from all machines. Furthermore, each machine is powerful enough to solve most optimization problems on their *local* dataset.



Notation: We will denote machine indices as superscripts (these are not powers!) and time indices as subscripts.

1.1 First attempt: independent optimization

One might be tempted to distribute the optimization problem by simply letting each machine compute the optimal solution on their local dataset, and then try to combine the solutions in some way. However, this approach in general has no hope without further assumptions about how the different datasets relate to each other. In the extreme, each solution might carry virtually zero information regarding the global optimum; for example, consider a multi-class classification problem where each machine only has access to training data from the same class. Then, a classifier that deterministically outputs the machine's class would be an optimal local solution and yet carry no meaning globally.

*These notes are class material that has not undergone formal peer review. The TAs and I are grateful for any reports of typos.

1.2 Second attempt: distributed gradient computation

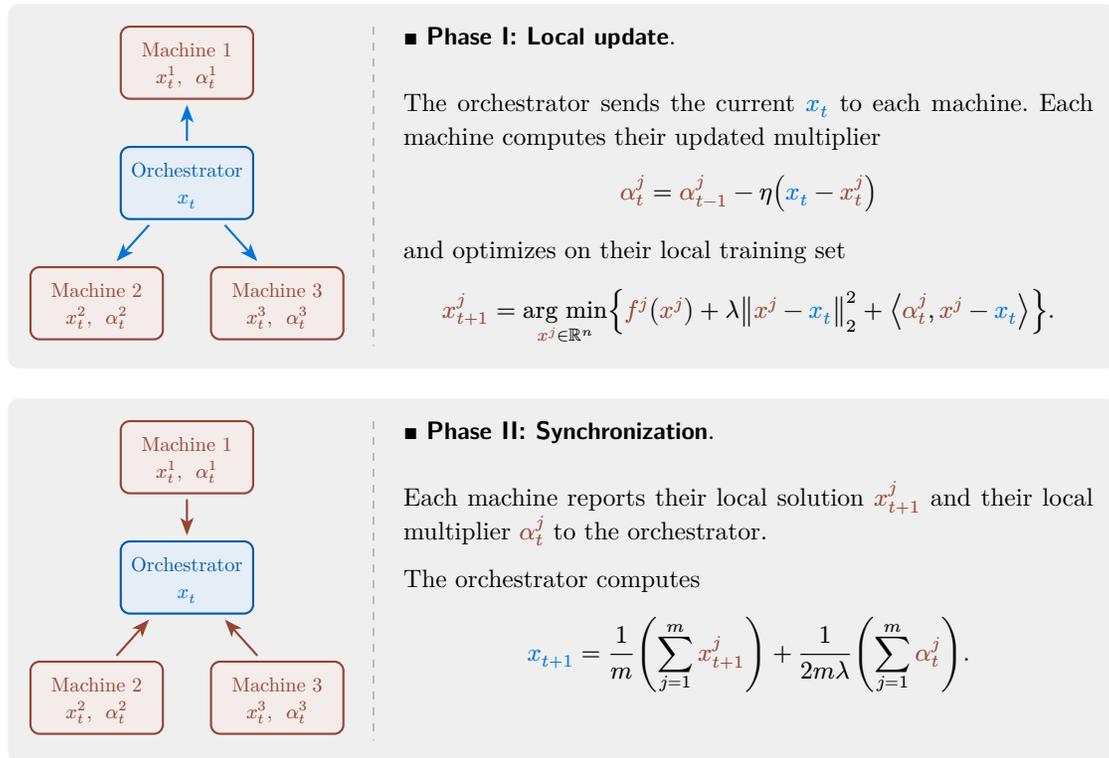
A more sophisticated approach is to compute the gradient of the loss function at each machine and then combine the gradients at a central location. This approach would work, but has two major drawbacks:

1. The total communication cost is mainly determined by the convergence rate of gradient descent. In particular, when the smoothness or the Lipschitzness of f is not very good, then the convergence rate of gradient descent is not very good. We will discuss today a solution that does not depend on the smoothness of f .
2. Gradient descent does not leverage the computation power of each machine very well. At each iteration, each machine simply computes a gradient. A more cost-effective approach would try to directly use the computational power of each machine to find some sort of optimal solution of the problem (while avoiding the naive approach of independent optimization discussed in Section 1.1).

2 The ADMM algorithm

The ADMM algorithm is a distributed optimization algorithm that combines the best of both worlds: it uses the computational power of each machine to find an optimal solution, while also ensuring that the global solution is meaningful.

The algorithm can be thought of as a two-step process parameterized by the learning rate $\eta > 0$ and a multiplier $\lambda > 0$, as follows.



2.1 ADMM as Lagrangian relaxation

At this point, one might wonder how the update rules of ADMM were picked. The answer is pretty ingenious. ADMM is based on the fundamental idea of rewriting the problem

$$v := \min_x \frac{1}{m} \sum_{j=1}^m f^j(x)$$

as the equivalent problem

$$\begin{aligned} \min_{x, x^j} \quad & \frac{1}{m} \sum_{j=1}^m \left(f^j(x^j) + \lambda \|x^j - x\|_2^2 \right) \\ \text{s.t.} \quad & x^j = x. \end{aligned} \tag{1}$$

The next step is to relax the equality constraint $x_{t+1}^j = x_{t+1}$ by introducing the *Lagrange relaxation* of the problem, in which the constraint is replaced by a penalization term in the objective function. In particular, the following key result shows that the maximum value of the Lagrange relaxation is the same as the value of the original problem when all f^j are convex.

Theorem 2.1. Assuming all f^j are convex, the optimal value of the original problem is equal to the maximum of the Lagrange function \mathcal{L} over all multipliers α^j , that is,

$$v = \max_{\alpha^1, \dots, \alpha^m} \min_{x, x^j} \left\{ \mathcal{L}(x, x^j; \alpha^j) := \frac{1}{m} \sum_{j=1}^m \left(f^j(x^j) + \lambda \|x^j - x\|_2^2 + \langle \alpha^j, x^j - x \rangle \right) \right\}.$$

In light of this result, the ADMM algorithm can be seen as a way to maximize the Lagrange function \mathcal{L} over the multipliers α^j by performing alternating optimization steps on different groups of variables:

- Assuming fixed x_{t-1} and x_{t-1}^j , the multipliers α_{t-1}^j are updated by gradient ascent, performing a step in the direction of the gradient of the objective function:

$$\alpha_t^j = \alpha_{t-1}^j + \eta (x_{t-1}^j - x_{t-1}) = \alpha_{t-1}^j - \eta (x_{t-1} - x_{t-1}^j).$$

- Assuming now fixed α_t^j and x_t , each machine now *solves* the optimization problem

$$x_{t+1}^j = \arg \min_{x^j \in \mathbb{R}^n} \left\{ f^j(x^j) + \lambda \|x^j - x_t\|_2^2 + \langle \alpha_t^j, x^j - x_t \rangle \right\}.$$

- Assuming now fixed α_t^j and x_{t+1}^j , the orchestrator computes the next iterate x_t by solving the optimization problem

$$\begin{aligned} x_{t+1} &= \arg \min_{x \in \mathbb{R}^n} \left\{ \frac{1}{m} \sum_{j=1}^m f^j(x^j) + \lambda \|x^j - x\|_2^2 + \langle \alpha_t^j, x^j - x \rangle \right\} \\ &= \arg \min_{x \in \mathbb{R}^n} \left\{ \frac{\lambda}{m} \sum_{j=1}^m \|x^j - x\|_2^2 - \langle \alpha_t^j, x \rangle \right\}. \end{aligned}$$

By setting the gradient with respect to x to 0 and solving for x , we obtain

$$\frac{1}{m} \sum_{j=1}^m (2\lambda(x - x^j) - \alpha_t^j) = 0$$

So, the optimal x_{t+1} is written in closed form as

$$x_{t+1} = \frac{1}{m} \sum_{j=1}^m x_{t+1}^j + \frac{1}{2m\lambda} \sum_{j=1}^m \alpha_t^j.$$

3 Analysis

The analysis of the correctness of ADMM can be broken down into several steps.

3.1 Part I: Orchestrating machine

Theorem 3.1. If the initial multipliers α_0^j are chosen such that $\sum_{j=1}^m \alpha_0^j = 0$, then, at all times t ,

$$\sum_{j=1}^m \alpha_t^j = 0, \quad \text{and} \quad x_{t+1} = \frac{1}{m} \sum_{j=1}^m x_{t+1}^j.$$

Proof. By definition,

$$\alpha_{t+1}^j = \alpha_t^j - \eta(x_{t+1} - x_{t+1}^j).$$

So, averaging over $j = 1, \dots, m$, we can write

$$\frac{1}{m} \sum_{j=1}^m \alpha_{t+1}^j = \frac{1}{m} \sum_{j=1}^m \alpha_t^j - \eta \left(x_{t+1} - \frac{1}{m} \sum_{j=1}^m x_{t+1}^j \right). \quad (2)$$

Using the definition of x_{t+1} , we also have

$$x_{t+1} = \left(\frac{1}{m} \sum_{j=1}^m x_{t+1}^j \right) + \frac{1}{2m\lambda} \sum_{j=1}^m \alpha_t^j \implies x_{t+1} - \frac{1}{m} \sum_{j=1}^m x_{t+1}^j = \frac{1}{2m\lambda} \sum_{j=1}^m \alpha_t^j. \quad (3)$$

Substituting the last expression into (2), we get

$$\frac{1}{m} \sum_{j=1}^m \alpha_{t+1}^j = \frac{1}{m} \sum_{j=1}^m \alpha_t^j - \frac{\eta}{2m\lambda} \sum_{j=1}^m \alpha_t^j = \frac{1}{m} \left(1 - \frac{\eta}{2\lambda} \right) \sum_{j=1}^m \alpha_t^j.$$

This shows that by induction, $\sum_{j=1}^m \alpha_t^j = 0$ at all t as long as the assumption holds for $t = 0$.

Plugging the expression $\sum_{j=1}^m \alpha_t^j = 0$ into (3), then yields

$$x_{t+1} = \frac{1}{m} \sum_{j=1}^m x_{t+1}^j$$

at all times, which is the second part of the claim. \square

3.2 Part II: Worker machines

Theorem 3.2 (Euclidean mirror descent lemma for ADMM). Let each $f^j : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex, and $\eta = 2\lambda$. Then, for any $y \in \mathbb{R}^n$ we have

$$\begin{aligned} \frac{1}{m} \sum_{j=1}^m f^j(x_{t+1}^j) &\leq f(y) + \frac{\eta}{2} \left(\|y - x_t\|_2^2 - \|y - x_{t+1}\|_2^2 - \|x_t - x_{t+1}\|_2^2 \right) \\ &\quad + \frac{1}{m} \sum_{j=1}^m \frac{1}{2\eta} \left(\|\alpha_t^j\|_2^2 - \|\alpha_{t+1}^j\|_2^2 - \|\alpha_{t+1}^j - \alpha_t^j\|_2^2 \right). \end{aligned}$$

Proof. Recall that each machine computes the update

$$x_{t+1}^j = \arg \min_{x^j \in \mathbb{R}^n} \left\{ f^j(x^j) + \lambda \|x^j - x_t\|_2^2 + \langle \alpha_t^j, x^j - x_t \rangle \right\}.$$

As the domain is open, the first-order necessary optimality conditions imply that

$$\nabla f^j(x_{t+1}^j) + 2\lambda(x_{t+1}^j - x_t) + \alpha_t^j = 0. \quad (4)$$

Since f^j is convex, the linear lower bound property of convex functions we can write

$$\begin{aligned} f^j(y) &\geq f^j(x_{t+1}^j) + \langle \nabla f^j(x_{t+1}^j), y - x_{t+1}^j \rangle \\ &= f^j(x_{t+1}^j) + \langle -2\lambda(x_{t+1}^j - x_t) - \alpha_t^j, y - x_{t+1}^j \rangle. \quad (\text{from (4)}) \end{aligned}$$

Rearranging the terms and using the hypothesis $\eta = 2\lambda$, we get

$$f^j(x_{t+1}^j) \leq f^j(y) + \eta \langle x_{t+1}^j - x_t, y - x_{t+1}^j \rangle + \langle \alpha_t^j, y - x_{t+1}^j \rangle. \quad (5)$$

At this point, we can use the update rule $\alpha_{t+1}^j = \alpha_t^j - \eta(x_{t+1} - x_{t+1}^j)$ to write

$$x_{t+1}^j = x_{t+1} + \frac{\alpha_{t+1}^j - \alpha_t^j}{\eta}.$$

Substituting the above expression into the first inner product of (5) then yields

$$f^j(x_{t+1}^j) \leq f^j(y) + \eta \langle x_{t+1} - x_t, y - x_{t+1}^j \rangle + \langle \alpha_{t+1}^j, y - x_{t+1}^j \rangle.$$

(Note that change in time index from α_t^j to α_{t+1}^j in the last term.) Substituting the expression for x_{t+1}^j into the last inner product, we finally obtain

$$f^j(x_{t+1}^j) \leq f^j(y) + \eta \langle x_{t+1} - x_t, y - x_{t+1}^j \rangle + \langle \alpha_{t+1}^j, y - x_{t+1} \rangle - \frac{1}{\eta} \langle \alpha_{t+1}^j, \alpha_{t+1}^j - \alpha_t^j \rangle. \quad (6)$$

From Theorem 3.1, we know that

$$\frac{1}{m} \sum_{j=1}^m \langle \alpha_{t+1}^j, y - x_{t+1} \rangle = 0, \quad \text{and} \quad \frac{1}{m} \sum_{j=1}^m \langle x_{t+1} - x_t, y - x_{t+1}^j \rangle = \langle x_{t+1} - x_t, y - x_{t+1} \rangle$$

for all t . Furthermore, we know that $f(y) = \frac{1}{m} \sum_{j=1}^m f^j(y)$. Hence, averaging (6) across all machines $j = 1, \dots, m$, we obtain

$$\frac{1}{m} \sum_{j=1}^m f^j(x_{t+1}^j) \leq f(y) + \eta \langle x_{t+1} - x_t, y - x_{t+1} \rangle - \frac{1}{m\eta} \sum_{j=1}^m \langle \alpha_{t+1}^j, \alpha_{t+1}^j - \alpha_t^j \rangle.$$

Substituting the three-point equalities

$$\begin{aligned} \langle x_{t+1} - x_t, y - x_{t+1} \rangle &= \frac{1}{2} \left(\|x_t - y\|_2^2 - \|x_{t+1} - y\|_2^2 - \|x_t - x_{t+1}\|_2^2 \right) \\ \text{and} \quad -\langle \alpha_{t+1}^j, \alpha_{t+1}^j - \alpha_t^j \rangle &= \frac{1}{2} \left(\|\alpha_t^j\|_2^2 - \|\alpha_{t+1}^j\|_2^2 - \|\alpha_{t+1}^j - \alpha_t^j\|_2^2 \right) \end{aligned}$$

yields the statement. \square

3.3 Putting the pieces together: telescoping step

The sum on the right-hand side of Theorem 3.2 telescopes. In particular, since $\alpha_0^j = 0$ for all j , we have

$$\frac{1}{mT} \sum_{t=0}^{T-1} \sum_{j=1}^m f^j(x_{t+1}^j) \leq f(y) + \frac{\eta}{2T} \|y - x_0\|_2^2 - \frac{1}{2\eta mT} \sum_{t=0}^{T-1} \sum_{j=1}^m \|\alpha_{t+1}^j - \alpha_t^j\|_2^2.$$

Since by definition $\alpha_{t+1}^j - \alpha_t^j = \eta(x_{t+1}^j - x_{t+1})$, the last term satisfies

$$-\frac{1}{2\eta mT} \sum_{t=0}^{T-1} \sum_{j=1}^m \|\alpha_{t+1}^j - \alpha_t^j\|_2^2 = -\frac{\eta}{2mT} \sum_{t=0}^{T-1} \sum_{j=1}^m \|x_{t+1}^j - x_{t+1}\|_2^2,$$

leading to

$$\frac{1}{mT} \sum_{t=0}^{T-1} \sum_{j=1}^m f^j(x_{t+1}^j) \leq f(y) + \frac{\eta}{2T} \|y - x_0\|_2^2 - \frac{\eta}{2mT} \sum_{t=0}^{T-1} \sum_{j=1}^m \|x_{t+1}^j - x_{t+1}\|_2^2.$$

Rearranging the terms and applying $y = x_*$, we therefore have

$$\frac{1}{T} \sum_{t=0}^{T-1} \left(\frac{1}{m} \sum_{j=1}^m f^j(x_{t+1}^j) + \frac{\eta}{2} \|x_{t+1}^j - x_{t+1}\|_2^2 \right) \leq f(x_*) + \frac{\eta}{2T} \|x_* - x_0\|_2^2.$$

Thus, we have the following guarantee.

Theorem 3.3. Let each $f^j : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex and nonnegative, $\lambda = \sqrt{T}$ and $\eta = 2\lambda$. Then, for any $T > 0$, at least one of the iterates $t \in \{0, \dots, T-1\}$ satisfies

$$\begin{aligned} \frac{1}{m} \sum_{j=1}^m f^j(x_{t+1}^j) &\leq f(x_*) + \frac{1}{\sqrt{T}} \|x_* - x_0\|_2^2, \\ \text{and} \quad \frac{1}{m} \sum_{j=1}^m \|x_{t+1}^j - x_{t+1}\|_2^2 &\leq \frac{1}{\sqrt{T}} f(x_*) + \frac{1}{T} \|x_* - x_0\|_2^2. \end{aligned}$$

So, as time progresses, the iterates x_{t+1}^j concentrate around the average x_{t+1} , and the average of the function values $f^j(x_{t+1}^j)$, which concentrates around $f(x_{t+1})$, converges to the optimal value at a rate of $\frac{1}{\sqrt{T}}$.

4 Further readings

The review article by [Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J., & others. \[Boy+11\]](#) is a modern and approachable introduction to ADMM and its variants and extensions.

[Boy+11] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, and others, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Foundations and Trends in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.