

Distributed Recovery/Regression/Classification using ADMM

By being very crafty with how we do the splitting, we can use ADMM to solve certain kinds of optimization programs in a distributed manner.

We consider (this material comes from [BPC⁺10, Sec. 8]) the general problem of “fitting” a vector $\mathbf{x} \in \mathbb{R}^N$ to an observed vector $\mathbf{b} \in \mathbb{R}^M$ through an $M \times N$ matrix \mathbf{A} . We will encourage \mathbf{x} to have certain structure using a regularizer. This type of problem is ubiquitous in signal processing and machine learning – the math stays the same, only the words change from area to area.

At a high level, we are interested in solving

$$\underset{\mathbf{x}}{\text{minimize}} \text{Loss}(\mathbf{A}\mathbf{x} - \mathbf{b}) + \text{Regularizer}(\mathbf{x})$$

where the $M \times N$ matrix \mathbf{A} and the vector \mathbf{b} are given. Notice that

$$\text{Loss}(\cdot) : \mathbb{R}^M \rightarrow \mathbb{R}, \quad \text{and} \quad \text{Regularizer}(\cdot) : \mathbb{R}^N \rightarrow \mathbb{R}.$$

We will assume that one or both of these functions are separable, at least at the block level. This means we can write

$$\text{Loss}(\mathbf{A}\mathbf{x} - \mathbf{b}) = \sum_{i=1}^B \ell_i(\mathbf{A}^{(i)}\mathbf{x} - \mathbf{b}^{(i)}),$$

where $\mathbf{A}^{(i)}$ are $M_i \times N$ matrices formed by partitioning the rows of \mathbf{A} , and $\mathbf{b}^{(i)} \in \mathbb{R}^{M_i}$ is the corresponding part of \mathbf{b} . For separable regularizers, we can write

$$\text{Regularizer}(\mathbf{x}) = \sum_{i=1}^C r_i(\mathbf{x}^{(i)}),$$

where the $\mathbf{x}^{(i)} \in \mathbb{R}^{N_i}$ partition the vector \mathbf{x} . These two types of separability will allow us to divide up the optimization in two different ways.

Example: Inverse Problems and Regression

Two popular methods for solving linear inverse problems and/or calculating regressors are solving

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

(*Tikhonov regularization* or *ridge regression*), and

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

(*the LASSO*).

These both clearly fit the separability criteria, as

$$\begin{aligned} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 &= \sum_{m=1}^M (\mathbf{a}_m^T \mathbf{x} - b[m])^2, \\ \|\mathbf{x}\|_2^2 &= \sum_{n=1}^N (x[n])^2 \\ \|\mathbf{x}\|_1 &= \sum_{n=1}^N |x[n]|. \end{aligned}$$

where \mathbf{a}_m^T is the m^{th} row of \mathbf{A} .

Example: Support Vector Machines

Previously, we saw how if we are given a set of M training examples (\mathbf{x}_m, y_m) , where $\mathbf{x}_m \in \mathbb{R}^N$ and $y_m \in \{-1, 1\}$, we can find a maximum margin linear classifier by solving

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{subject to} \quad y_m(b - \langle \mathbf{x}_m, \mathbf{w} \rangle) + 1 \leq 0, \quad m = 1, \dots, M.$$

With the classifier trained (optimal solution \mathbf{w}^*, b^* computed), we can assign a label y' to a new point \mathbf{x}' using

$$y' = \text{sign}(\langle \mathbf{x}', \mathbf{w}^* \rangle + b^*).$$

Instead of enforcing the constraints above strictly, we can allow some errors by penalizing mis-classifications on the training data appropriately. One reasonable way to do this is make the loss zero if $y_m(b - \langle \mathbf{x}_m, \mathbf{w} \rangle) + 1 \leq 0$, and then have it increase linearly as this quantity exceeds zero. That is, we solve

$$\min_{\mathbf{w}, b} \sum_{m=1}^M \ell(y_m(b - \langle \mathbf{x}_m, \mathbf{w} \rangle) + 1) + \frac{1}{2} \|\mathbf{w}\|_2^2,$$

where $\ell(\cdot)$ is

$$\ell(u) = (u)_+ = \begin{cases} 0, & u \leq 0, \\ u, & u > 0. \end{cases}$$

This is penalty is often called the **hinge loss**. Note that the argument for $\ell(\cdot)$ is an affine function of the optimization variables:

$$y_m(b - \langle \mathbf{x}_m, \mathbf{w} \rangle) + 1 = \begin{bmatrix} -y_m \mathbf{x}_m^T & y_m \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} + 1.$$

Both the loss function and regularizer in this formulation of the SVM are clearly separable.

Splitting across examples

This framework is useful when we have “many measurements of a small vector” or “large volumes of low-dimensional data”.

We partition the rows of \mathbf{A} and entries of \mathbf{b} :

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}^{(1)} \\ \mathbf{A}^{(2)} \\ \vdots \\ \mathbf{A}^{(B)} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \mathbf{b}^{(1)} \\ \mathbf{b}^{(2)} \\ \vdots \\ \mathbf{b}^{(B)} \end{bmatrix}.$$

If the loss function is separable over this partition, our optimization problem is

$$\underset{\mathbf{x}}{\text{minimize}} \quad \sum_{i=1}^B \ell_i(\mathbf{A}^{(i)} \mathbf{x} - \mathbf{b}^{(i)}) + r(\mathbf{x}),$$

where $r(\cdot)$ is the regularizer. We start by splitting the optimization variables in the loss function and those in the regularizer, arriving at the equivalent program

$$\underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} \quad \sum_{i=1}^B \ell_i(\mathbf{A}^{(i)} \mathbf{x} - \mathbf{b}^{(i)}) + r(\mathbf{z}) \quad \text{subject to} \quad \mathbf{x} - \mathbf{z} = \mathbf{0}.$$

This does not make the Lagrangian for the primal update separable, as the \mathbf{A}_i are still tying together all of the entries in \mathbf{x} . The trick is to introduce B different vectors $\mathbf{x}^{(i)} \in \mathbb{R}^N$, one for each block, and then use the constraints to make them all agree. This is done with

$$\begin{aligned} & \underset{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(B)}, \mathbf{z}}{\text{minimize}} \quad \sum_{i=1}^B \ell_i(\mathbf{A}^{(i)} \mathbf{x}^{(i)} - \mathbf{b}^{(i)}) + r(\mathbf{z}) \\ & \text{subject to} \quad \mathbf{x}^{(i)} - \mathbf{z} = \mathbf{0}, \quad i = 1, \dots, B. \end{aligned}$$

The augmented Lagrangian for this last problem can be expressed as

$$\mathcal{L}_\rho(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(B)}, \mathbf{z}, \boldsymbol{\mu}^{(1)}, \dots, \boldsymbol{\mu}^{(B)}) = \sum_{i=1}^B \mathcal{L}_i(\mathbf{x}^{(i)}, \mathbf{z}, \boldsymbol{\mu}^{(i)}),$$

where

$$\mathcal{L}_i(\mathbf{x}^{(i)}, \mathbf{z}, \boldsymbol{\mu}^{(i)}) = \ell_i(\mathbf{A}^{(i)} \mathbf{x}^{(i)} - \mathbf{b}^{(i)}) + \frac{r(\mathbf{z})}{B} + \frac{\rho}{2} \|\mathbf{x}^{(i)} - \mathbf{z} + \boldsymbol{\mu}^{(i)}\|_2^2$$

and the $\boldsymbol{\mu}^{(i)}$ are the (rescaled) Lagrange multipliers for the constraint $\mathbf{x}^{(i)} - \mathbf{z} = \mathbf{0}$.

As the Lagrangian is separable over the B blocks, each of the primal updates for the \mathbf{x}_i can be performed independently. This makes the ADMM iteration

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

The \mathbf{z} update can be written in terms of the average of the $\mathbf{x}_{k+1}^{(i)}$ and the $\boldsymbol{\mu}_k^{(i)}$. To see this, first note that

$$\begin{aligned} \sum_{i=1}^B \|\mathbf{z} - \mathbf{v}_i\|_2^2 &= B\|\mathbf{z}\|_2^2 - 2 \left\langle \mathbf{z}, \sum_{i=1}^B \mathbf{v}_i \right\rangle + \sum_{i=1}^B \|\mathbf{v}_i\|_2^2 \\ &= B\|\mathbf{z}\|_2^2 - 2B \langle \mathbf{z}, \bar{\mathbf{v}} \rangle + B\|\bar{\mathbf{v}}\|_2^2 + \left(-B\|\bar{\mathbf{v}}\|_2^2 + \sum_{i=1}^B \|\mathbf{v}_i\|_2^2 \right) \\ &= B\|\mathbf{z} - \bar{\mathbf{v}}\|_2^2 + \left(-B\|\bar{\mathbf{v}}\|_2^2 + \sum_{i=1}^B \|\mathbf{v}_i\|_2^2 \right). \end{aligned}$$

where $\bar{\mathbf{v}} = \frac{1}{B} \sum_{i=1}^B \mathbf{v}_i$. Thus

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

Distributed ADMM (dividing rows of \mathbf{A})

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

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

$$\boldsymbol{\mu}_{k+1}^{(i)} = \boldsymbol{\mu}_k^{(i)} + \mathbf{x}_{k+1}^{(i)} - \mathbf{z}_{k+1}$$

where

$$\bar{\mathbf{x}}_{k+1} = \frac{1}{B} \sum_{i=1}^B \mathbf{x}_{k+1}^{(i)}, \quad \bar{\boldsymbol{\mu}}_k = \frac{1}{B} \sum_{i=1}^B \boldsymbol{\mu}_k^{(i)}.$$

The high-level architecture is that B separate units solve independent optimization programs for the B $\mathbf{x}^{(i)}$ updates. These are collected and averaged, and a single optimization program is solved to get the \mathbf{z} update. The new \mathbf{z} is then communicated back to each of the B units. The Lagrange multiplier update can be easily computed both centrally and at the B units, so these do not have to be communicated.

Example: The LASSO

With $\ell_i(\mathbf{A}^{(i)}\mathbf{x}^{(i)} - \mathbf{b}^{(i)}) = \|\mathbf{A}^{(i)}\mathbf{x}^{(i)} - \mathbf{b}^{(i)}\|_2^2$ and $r(\mathbf{z}) = \tau\|\mathbf{z}\|_1$, the ADMM iteration becomes

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

The $\mathbf{x}^{(i)}$ updates are all small unconstrained least-squares problems whose solutions can be computed independently; the \mathbf{z} update is a simple soft thresholding, and the $\boldsymbol{\mu}^{(i)}$ and $\bar{\boldsymbol{\mu}}$ updates are computed simply by adding vectors.

Example: SVMs

For the SVM, we collect the weights and the offset into a single optimization vector

$$\mathbf{x} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} \in \mathbb{R}^{N+1}$$

and set

$$\mathbf{A} = \begin{bmatrix} -y_1\mathbf{x}_1^T & y_1 \\ \vdots & \vdots \\ -y_M\mathbf{x}_M^T & y_M \end{bmatrix}$$

If we partition the data (\mathbf{A}) into B blocks ($\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(B)}$) then we can express the i^{th} component of the augmented Lagrangian as

$$\mathcal{L}_i(\mathbf{x}^{(i)}, \mathbf{z}, \boldsymbol{\mu}^{(i)}) = \mathbf{1}^T(\mathbf{A}^{(i)}\mathbf{x}^{(i)} + \mathbf{1})_+ + \frac{r(\mathbf{z})}{B} + \frac{\rho}{2} \|\mathbf{x}^{(i)} - \mathbf{z} + \boldsymbol{\mu}^{(i)}\|_2^2.$$

Note that the regularization does not include the last term in \mathbf{z} :

$$r(\mathbf{z}) = \frac{1}{2} \sum_{n=1}^N |z[n]|^2$$

This results in the ADMM iteration

$$\mathbf{x}_{k+1}^{(i)} = \arg \min_{\mathbf{x}^{(i)}} \left(\mathbf{1}^\top (\mathbf{A}^{(i)} \mathbf{x}^{(i)} + \mathbf{1})_+ + \frac{\rho}{2} \|\mathbf{x}^{(i)} - \mathbf{z}_k + \boldsymbol{\mu}_k^{(i)}\|_2^2 \right),$$
$$\mathbf{z}_{k+1}[n] = \begin{cases} \frac{B\rho}{1+B\rho} (\bar{\mathbf{x}}_{k+1}[n] + \bar{\boldsymbol{\mu}}_k[n]), & n = 1, \dots, N, \\ \bar{\mathbf{x}}_{k+1}[n] + \bar{\boldsymbol{\mu}}_k[n], & n = N + 1, \end{cases}$$
$$\boldsymbol{\mu}_{k+1}^{(i)} = \boldsymbol{\mu}_k^{(i)} + \mathbf{x}_{k+1}^{(i)} - \mathbf{z}_{k+1}.$$

Splitting across features

Similarly, we can divide up the *columns* of \mathbf{A} . This is described in [BPC⁺10, Section 8.3].

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.