

ADMM in PyTorch

Milan Parmar (IIT2021121), Sibasish Behera (IIT2021136), Bittu Sharma (IIT2021200),
Kushagra Jain (IIT2021191)

Indian Institute of Information Technology Allahabad

14, Nov 2024

Abstract

The Alternating Direction Method of Multipliers (ADMM) is an optimization algorithm which solves the objective function by breaking it into smaller pieces. This allows us to pass the smaller pieces to separate machines/processors and solve them parallelly. We implemented the ADMM serialized and parallel code for Lasso and Ridge objective functions in PyTorch framework.

1 Introduction

Alternating Direction Method of Multipliers [1] is a simple but powerful algorithm (framework) which works really well for distributed convex optimization. The core idea behind it is broadcasting and gathering, where broadcasting means splitting the objective function into smaller parts and giving it to multiple processors/machines for solving, and gathering means gathering the results from these processors/machines and combining them for getting the update rule for the next iteration. The advantage of this algorithm is that there are no constraints on the objective function other than convexity. We do not even need the loss function to be differentiable.

The flow of the report is as follows: We first introduce the existing state of the art implementations for ADMM. Followed by that, we explain the working of ADMM and then move on to our contribution, which is implementation of ADMM for Lasso and Ridge in PyTorch (serial and parallel) and their comparison with sklearn's implementation. Finally, we conclude our report and state our ongoing work and future scope.

2 ADMM Algorithm

The ADMM objective is solving a sum of functions with linear constraints. In general, for parallelizing the minimization process, we divide the sum of the objective function into N parts as follows:

$$\begin{aligned} & \min_{x(i), Z} \sum_{i=1}^N f_i(x(i)) \\ \text{st} \quad & Z = x(i) \end{aligned}$$

The algorithm consists of updating three main steps:

1. Minimization of x

$$x_{t+1}^{(i)} = \underset{x}{\operatorname{argmin}} f_i(x) + \nu_t^{(i)}(x^{(i)} - z_t) + \frac{\rho}{2}\|x^{(i)} - z_t\|_2^2$$

2. Minimization of z

$$z_{t+1}^{(i)} = \frac{1}{N} \sum_{i=1}^N (x_{t+1}^{(i)} + \frac{1}{\rho} \nu_t^{(i)})$$

$$z_{t+1} = \bar{x}_{t+1} + \frac{1}{\rho} \bar{\nu}_t$$

$$\text{where, } \bar{x}_{t+1} = \frac{1}{N} \sum_{i=1}^N x_{t+1}^{(i)} \text{ and } \bar{\nu}_{t+1} = \frac{1}{N} \sum_{i=1}^N \nu_{t+1}^{(i)}$$

3. Update of dual variable ν_t

$$\nu_{t+1}^{(i)} = \nu_t^{(i)} + \rho(x_{t+1}^{(i)} - z_{t+1}^{(i)})$$

$$\bar{\nu}_{t+1} = \bar{\nu}_t + \rho(\bar{x}_{t+1} - \bar{z}_{t+1})$$

Now, the updates for each $x_{t+1}^{(i)}$ and $\nu_{t+1}^{(i)}$ can be done in parallel

3 LASSO Regression

The LASSO objective function is a modified Linear Regression objective function. It involves an L1 regularization term. So the modified objective is:

$$\min_x \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$$

where $x \in \mathbb{R}^d$ is a vector of weights of the model that need to be tuned and $A \in \mathbb{R}^{n \times d}$ is a matrix of all the input features and $b \in \mathbb{R}^n$ is the vector of all outputs.

The above function can be written in multiple ways as an ADMM objective function subject to linear constraints. One way to do that is the Serial Vectorized implementation which is as follows:

3.1 Vectorized Serial Implementation

We can split the objective function by renaming one of the two variables and introduce an equality constraint to bring in the ADMM form. The modified LASSO objective is:

$$\min_{x, z} \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|z\|_1$$

$$\text{st } x = z$$

The update rule for each of the two variables is as follows:

1. Minimization of x

$$x_{t+1} = \underset{x}{\operatorname{argmin}} \quad \frac{1}{2} \|Ax - b\|_2^2 + \frac{\rho}{2} \|x - z_t + \frac{1}{\rho} \nu_t\|_2^2$$

which has a closed form solution

$$x_{t+1} = (A^T A + \rho I)^{-1} (A^T b + \rho z_t - \nu_t)$$

2. Minimization of z

$$z_{t+1} = \underset{z}{\operatorname{argmin}} \quad \lambda \|z\|_1 + \frac{\rho}{2} \|x_{t+1} - z + \frac{1}{\rho} \nu_t\|_2^2$$

$$z_{t+1} = x_{t+1} + \frac{\nu_t}{\rho} - \frac{\lambda}{\rho} \operatorname{sign}(z)$$

3. Update of dual variable ν_t

$$\nu_{t+1} = \nu_t + \rho(x_{t+1} - z_{t+1})$$

3.2 Parallel Implementation

We can rewrite the objective function of LASSO for running in distributed/parallel fashion by splitting across all the samples:

$$\min_{x^{(i)}, z} \quad \frac{1}{2} \sum_{i=1}^N \|A^{(i)} x^{(i)} - b^{(i)}\|_2^2 + \lambda \|z\|_1$$

$$\text{st} \quad x^{(i)} - z = 0, \text{ for } i = 1, \dots, N$$

The modified update rule for each variable is as follows:

1. Minimization of x_i

$$x_{t+1}^{(i)} = \underset{x^{(i)}}{\operatorname{argmin}} \quad \frac{1}{2} \|A^{(i)} x^{(i)} - b^{(i)}\|_2^2 + \frac{\rho}{2} \|x^{(i)} - z_t + \frac{1}{\rho} \nu_t^{(i)}\|_2^2$$

which has a closed form solution

$$x_{t+1}^{(i)} = (A^{(i)T} A^{(i)} + \rho I)^{-1} (A^{(i)T} b^{(i)} + \rho z_t - \nu_t^{(i)}), \text{ for } i = 1, \dots, N$$

2. Minimization of z

$$z_{t+1} = \underset{z}{\operatorname{argmin}} \quad \lambda \|z\|_1 + \frac{\rho}{2} \|\bar{x}_{t+1} - z + \frac{1}{\rho} \bar{\nu}_t\|_2^2$$

$$z_{t+1} = \bar{x}_{t+1} + \frac{\bar{\nu}_t}{\rho} - \frac{\lambda}{\rho} \operatorname{sign}(z)$$

3. Update of dual variable ν_t

$$\nu_{t+1}^{(i)} = \nu_t^{(i)} + \rho(x_{t+1}^{(i)} - z_{t+1})$$

4 Ridge Regression

The Ridge Regression objective is similar to the LASSO objective, but in this case we use the L-2 norm squared instead of L-1 norm as regularization term with λ . It is as follows:

$$\min_x \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_2^2$$

where $x \in \mathbb{R}^d$ is a vector of weights of the model that need to be tuned and $A \in \mathbb{R}^{n \times d}$ is a matrix of all the input features and $b \in \mathbb{R}^n$ is the vector of all outputs.

Similar to LASSO we implemented this algorithm in the following two ways:

4.1 Vectorized Serial Implementation

We can split the objective function by renaming one of the two variables and introduce an equality constraint to bring in the ADMM form. The modified Ridge objective is:

$$\begin{aligned} \min_{x, z} \quad & \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|z\|_2^2 \\ \text{st} \quad & x = z \end{aligned}$$

The update rule for each of the two variables is as follows:

1. Minimization of x

$$x_{t+1} = \underset{x}{\operatorname{argmin}} \quad \frac{1}{2} \|Ax - b\|_2^2 + \frac{\rho}{2} \|x - z_t + \frac{1}{\rho} \nu_t\|_2^2$$

which has a closed form solution

$$x_{t+1} = (A^T A + \rho I)^{-1} (A^T b + \rho z_t - \nu_t)$$

2. Minimization of z

$$\begin{aligned} z_{t+1} = \underset{z}{\operatorname{argmin}} \quad & \lambda \|z\|_2^2 + \frac{\rho}{2} \|x_{t+1} - z + \frac{1}{\rho} \nu_t\|_2^2 \\ z_{t+1} = \quad & \frac{\rho x_{t+1} + \nu_t}{2\lambda + \rho} \end{aligned}$$

3. Update of dual variable ν_t

$$\nu_{t+1} = \nu_t + \rho(x_{t+1} - z_{t+1})$$

4.2 Parallel Implementation

We can rewrite the objective function of Ridge for running in distributed/parallel fashion by splitting across all the samples:

$$\begin{aligned} \min_{x^{(i)}, z} \quad & \frac{1}{2} \sum_{i=1}^N \|A^{(i)}x^{(i)} - b^{(i)}\|_2^2 + \lambda \|z\|_2^2 \\ \text{st} \quad & x^{(i)} - z = 0, \text{ for } i = 1, \dots, N \end{aligned}$$

The modified update rule for each of the two variables is as follows:

1. Minimization of x_i

$$x_{t+1}^{(i)} = \underset{x^{(i)}}{\operatorname{argmin}} \quad \frac{1}{2} \|A^{(i)}x^{(i)} - b^{(i)}\|_2^2 + \frac{\rho}{2} \|x^{(i)} - z_t + \frac{1}{\rho} \nu_t^{(i)}\|_2^2$$

which has a closed form solution

$$x_{t+1}^{(i)} = (A^{(i)T}A^{(i)} + \rho I)^{-1}(A^{(i)T}b^{(i)} + \rho z_t - \nu_t^{(i)}), \text{ for } i = 1, \dots, N$$

2. Minimization of z

$$\begin{aligned} z_{t+1} &= \underset{z}{\operatorname{argmin}} \quad \lambda \|z\|_2^2 + \frac{\rho}{2} \|\bar{x}_{t+1} - z + \frac{1}{\rho} \bar{\nu}_t\|_2^2 \\ z_{t+1} &= \frac{\rho \bar{x}_{t+1} + \bar{\nu}_t}{2\lambda + \rho} \end{aligned}$$

3. Update of dual variable ν_t

$$\nu_{t+1}^{(i)} = \nu_t^{(i)} + \rho(x_{t+1}^{(i)} - z_{t+1})$$

5 Results

We tested our parallel implementation of ADMM for LASSO and Ridge against the sklearn's implementation on the Diabetes dataset [4].

5.1 Diabetes Dataset

- The diabetes dataset has 10 features:
age, sex, body mass index, average blood pressure, and six blood serum measurements (S1-S6)
- The output is a quantitative measure of disease progression one year after the measurements of the 10 features
- Number of samples: 442

5.2 LASSO Regression

We ran our ADMM Lasso solver on above Diabetes dataset: Figure 1(a) shows objective values of ADMM Lasso solver.

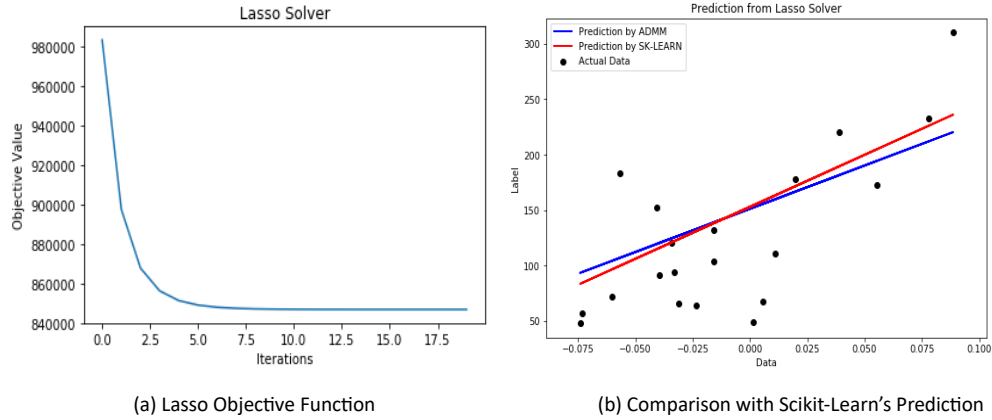


Figure 1: Lasso Solver

Figure 1(b) shows prediction comparison with Sklearn's lasso solver.

5.3 Ridge Regression

We ran our ADMM Ridge solver on above Diabetes dataset: Figure 2(a) shows objective values of ADMM Ridge solver and figure 2(b) shows prediction comparison with Sklearn's Ridge solver.

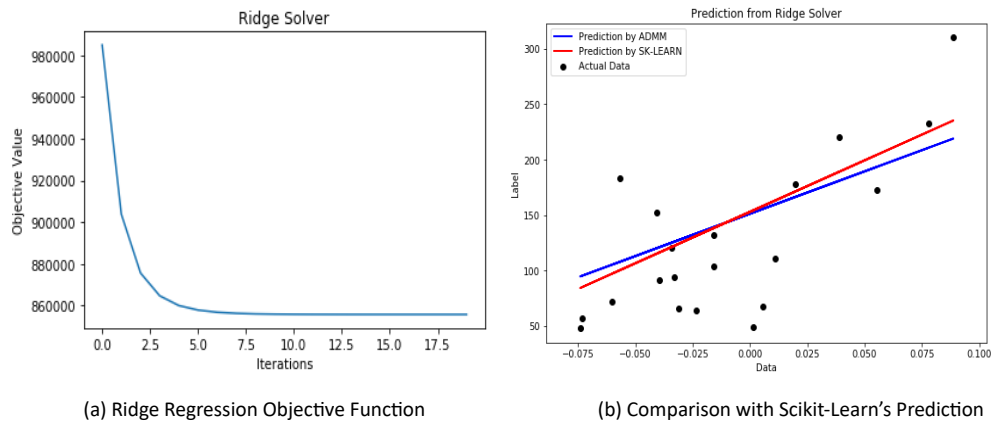


Figure 2: Ridge Solver

The following plot shows a comparison of how our ADMM performs against gradient descent and newton algorithms for a Ridge objective function.

6 Conclusion

We implemented the above mentioned two algorithms in the PyTorch framework. There are many such algorithms that can be optimized with ADMM. Our implementation gives similar results to the state of the art solvers from scikit learn. As expected, the ADMM implementation performs faster than the scikit learn implementation as seen from the results.

References

- [1] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein (2011) Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers. *Foundations and Trends in Machine Learning* 3(1):1–122.
- [2] Stanford University Convex Optimization Group, cvxflow, <https://github.com/cvxgrp/cvxflow/blob/master/cvxflow/solvers/admm.py>, 2017