

Training Neural Networks with ADMM

Project 3 Report

Team 6

CHAI DI

20458239

Big Data Technology

HKUST

dchai@connect.ust.hk

LU GUANNAN

20454477

Big Data Technology

HKUST

gluab@connect.ust.hk

FENG YAJUN

20404460

Big Data Technology

HKUST

yfengaj@connect.ust.hk

LIN YUWEI

20446054

Big Data Technology

HKUST

ylinbq@connect.ust.hk

1. Problem Description

Nowadays, training neural networks always accompanies with big data volume and large number of iterations, which places a greater demand on efficiency and effectiveness of computation. However, rather than blindly using GPUs and increasing computation speed, we can trace back to the source and figure out some innovative gradient descent methods to solve the optimization problems. Actually, common-used stochastic gradient descent method has its drawbacks: all gradient-based approaches suffer from vanishing gradient problems and easily get stuck at local minima and saddle points. So in the paper “Training Neural Networks without Gradients: A Scalable ADMM Approach”, it separates the objective function at each layer of a neural network into two terms that can be updated and solved as a set of minimization sub-problems without gradients, which to some extent helps alleviate the weaknesses met by stochastic gradient descent. This project aims to implement the optimization method in above paper. By training a deep neural network, we compare ADMM (Alternating Direction Methods of Multipliers) and standard back propagation method on image datasets. In this report, section 1 is the introduction of the project. Section 2 introduces the mathematical principles and algorithmic details of ADMM. Section 3 focus on the experiment process and results. Section 4 is the comparison and conclusion.

2. Principles

Given input vector \mathbf{a}_{l-1} , linear weight vector \mathbf{W}_l and a non-linear neural activation function \mathbf{h}_l , the output is $\mathbf{a}_l = \mathbf{h}_l(\mathbf{W}_l \mathbf{a}_{l-1})$. A neural network can be represented as

$$\mathbf{f}(\mathbf{a}_0; \mathbf{W}) = \mathbf{W}_k \mathbf{h}_{k-1}(\mathbf{W}_{k-1} \mathbf{h}_{k-2}(\dots(\mathbf{W}_2 \mathbf{h}_1(\mathbf{W}_1 \mathbf{a}_0))))$$

where $\mathbf{W} = \{\mathbf{W}_l\}$ is the weight matrix.

The optimization objective function is

$$\min_{\mathbf{W}} l(\mathbf{f}(\mathbf{a}_0; \mathbf{W}), \mathbf{y})$$

To decouple the weights \mathbf{W} from nonlinear functions \mathbf{h}_l , we get

$$\begin{aligned} \min_{\{\mathbf{W}_l\}, \{\mathbf{a}_l\}, \{\mathbf{z}_l\}} \quad & l(\mathbf{z}_L, \mathbf{y}) \\ \text{s. t.} \quad & \mathbf{z}_l = \mathbf{W}_l \mathbf{a}_{l-1}, \quad l = 1, 2, \dots, L \\ & \mathbf{a}_l = \mathbf{h}_l(\mathbf{z}_l), \quad l = 1, 2, \dots, L \end{aligned}$$

The Augmented Lagrangian function (adding a \mathbf{l}_2 penalty) is

$$\begin{aligned} & \min_{\{\mathbf{W}_l\}, \{\mathbf{a}_l\}, \{\mathbf{z}_l\}} L_\lambda(\mathbf{W}, \mathbf{a}, \mathbf{z}) \\ = & \min_{\{\mathbf{W}_l\}, \{\mathbf{a}_l\}, \{\mathbf{z}_l\}} l(\mathbf{z}_L, \mathbf{y}) + \langle \mathbf{z}_L, \boldsymbol{\lambda} \rangle + \boldsymbol{\beta}_L \|\mathbf{z}_L - \mathbf{W}_L \mathbf{a}_{L-1}\|_2^2 \\ & + \sum_{l=1}^{L-1} [\gamma_l \|\mathbf{a}_l - \mathbf{h}_l(\mathbf{z}_l)\|_2^2 + \boldsymbol{\beta}_l \|\mathbf{z}_l - \mathbf{W}_l \mathbf{a}_{l-1}\|_2^2] \end{aligned}$$

Where $\boldsymbol{\lambda}$ is a vector of Lagrange multipliers.

The augmented Lagrangian algorithm employing an alternating direction gives

$$\begin{cases} \mathbf{W}_l = \underset{\mathbf{W}_l}{\operatorname{argmin}} L_\lambda(\mathbf{W}_l, \mathbf{a}_l, \mathbf{z}_l) \\ \mathbf{a}_l = \underset{\mathbf{a}_l}{\operatorname{argmin}} L_\lambda(\mathbf{W}_l, \mathbf{a}_l, \mathbf{z}_l) \\ \mathbf{z}_l = \underset{\mathbf{z}_l}{\operatorname{argmin}} L_\lambda(\mathbf{W}_l, \mathbf{a}_l, \mathbf{z}_l) \\ \boldsymbol{\lambda} = \boldsymbol{\lambda} + \boldsymbol{\beta}_L \mathbf{z}_L - \mathbf{W}_L \mathbf{a}_{L-1} \end{cases}$$

To respectively solve above update iterations,

(1) \mathbf{W}_l

$$\begin{aligned} \mathbf{W}_l &= \underset{\mathbf{W}_l}{\operatorname{argmin}} \boldsymbol{\beta}_l \|\mathbf{z}_l - \mathbf{W}_l \mathbf{a}_{l-1}\|_2^2 \\ \mathbf{W}_l &\leftarrow \mathbf{z}_l \mathbf{a}_{l-1}^+ \end{aligned}$$

(2) \mathbf{a}_l

$$\mathbf{a}_l = \underset{\mathbf{a}_l}{\operatorname{argmin}} \gamma_l \|\mathbf{a}_l - \mathbf{h}_l(\mathbf{z}_l)\|_2^2 + \beta_{l+1} \|\mathbf{z}_{l+1} - \mathbf{W}_{l+1} \mathbf{a}_l\|_2^2$$

Taking gradient and setting it to 0, we get

$$(\beta_{l+1} \mathbf{W}_{l+1}^T \mathbf{W}_{l+1} + \gamma_l) \mathbf{a}_l = \beta_{l+1} \mathbf{W}_{l+1}^T \mathbf{z}_{l+1} + \gamma_l \mathbf{h}_l(\mathbf{z}_l),$$

which leads to

$$\mathbf{a}_l \leftarrow (\beta_{l+1} \mathbf{W}_{l+1}^T \mathbf{W}_{l+1} + \gamma_l)^{-1} (\beta_{l+1} \mathbf{W}_{l+1}^T \mathbf{z}_{l+1} + \gamma_l \mathbf{h}_l(\mathbf{z}_l))$$

(3) \mathbf{z}_l

$$\mathbf{z}_l \leftarrow \underset{\mathbf{z}_l}{\operatorname{argmin}} \gamma_l \|\mathbf{a}_l - \mathbf{h}_l(\mathbf{z}_l)\|_2^2 + \beta_l \|\mathbf{z}_l - \mathbf{W}_l \mathbf{a}_{l-1}\|_2^2$$

which is a non-convex function but can fortunately be solved in closed form.

(4) λ

After updating $\{\mathbf{W}_l\}, \{\mathbf{a}_l\}, \{\mathbf{z}_l\}$, the Lagrange multiplier update is given by

$$\lambda \leftarrow \lambda + \beta_L \mathbf{z}_L - \mathbf{W}_L \mathbf{a}_{L-1}$$

So the algorithm of ADMM is shown below:

Algorithm 1 ADMM for Neural Nets

Input: training features $\{\mathbf{a}_0\}$, and labels $\{y\}$,
Initialize: allocate $\{\mathbf{a}_l\}_{l=1}^{L-1}$, $\{\mathbf{z}_l\}_{l=1}^L$, and λ
repeat
 for $l = 1, 2, \dots, L - 1$ **do**
 $\mathbf{W}_l \leftarrow \mathbf{z}_l \mathbf{a}_{l-1}^\dagger$
 $\mathbf{a}_l \leftarrow (\beta_{l+1} \mathbf{W}_{l+1}^T \mathbf{W}_{l+1} + \gamma_l I)^{-1} (\beta_{l+1} \mathbf{W}_{l+1}^T \mathbf{z}_{l+1} + \gamma_l \mathbf{h}_l(\mathbf{z}_l))$
 $\mathbf{z}_l \leftarrow \arg \min_z \gamma_l \|\mathbf{a}_l - \mathbf{h}_l(\mathbf{z})\|_2^2 + \beta_l \|\mathbf{z}_l - \mathbf{W}_l \mathbf{a}_{l-1}\|_2^2$
 end for
 $\mathbf{W}_L \leftarrow \mathbf{z}_L \mathbf{a}_{L-1}^\dagger$
 $\mathbf{z}_L \leftarrow \arg \min_z \ell(\mathbf{z}, y) + \langle \mathbf{z}_L, \lambda \rangle + \beta_L \|\mathbf{z}_L - \mathbf{W}_L \mathbf{a}_{L-1}\|_2^2$
 $\lambda \leftarrow \lambda + \beta_L (\mathbf{z}_L - \mathbf{W}_L \mathbf{a}_{L-1})$
until converged

3. Experiments

3.1 Datasets

The dataset used in this project, is from Yann LeCun, Corinna Cortes, Christopher J.C. Burges: “THE MNIST DATABASE”. It is a set of labeled hand-written digits, including a

training set of 60000 datasets in the size of 28*28 pixels, along with a training set labels. And this project utilizes the datasets labeled 0 or 1, which is 12665 out of 60000. Among the 12665 original datasets, 10000 of them are used as training set in this project, and the remaining 2665 sets are for testing.

3.2 Process

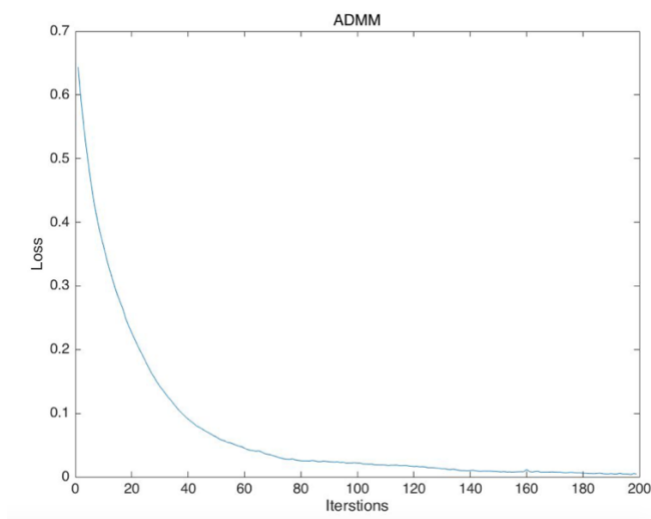
Firstly, a simple preprocessing is conducted. The original dataset is 60000 28*28 images, after filtering labels with 0 or 1, 12665 28*28 images remain, then the 12665 matrix are reshaped to a vector with the length of 784.

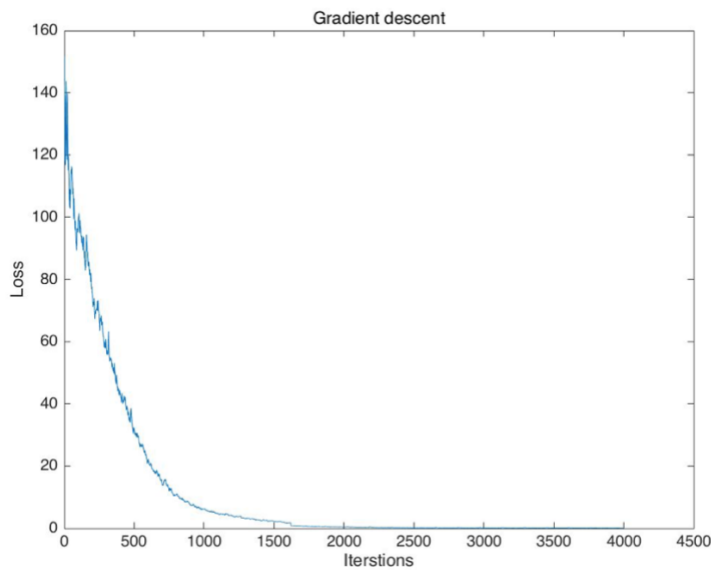
Afterwards, a fully connected neural network is built up, with size 100, 50, 10 and 1 sequentially, along with the mean square error function to calculate the loss. In the network, different parameter update functions, which are ADMM and Gradient descent, are implemented. To compare the performance, the evaluation criteria accuracy and error rate are calculated in each iterations.

Finally, the result in the accuracy aspect is 0.99061913696 for ADMM compared to 0.994746716689 for Gradient descent.

4. Comparison

The charts below show a comparison of performance between ADMM and Gradient descent in the form of loss.





1. The charts show that the ADMM is much smoother compared with Gradient descent. In the whole experiment process, for Gradient descent, the accuracy is quite different in each run, meanwhile there are occasions when the loss is decreasing while the accuracy remains.
2. Compared with Gradient, the loss with ADMM decreases much faster. Considering loss, ADMM reaches below 0.5 after about 60 iterations, while Gradient descent needs about 1500 iterations to get the same result.

5. Reference

5.1 Yann LeCun, Corinna Cortes, Christopher J.C. Burges: "THE MNIST DATABASE".

<http://yann.lecun.com/exdb/mnist/>

5.2 Gavin Taylor, Ryan Burmeister, Zheng Xu, Bharat Singh, Ankit Patel, Tom Goldstein: Proceedings of The 33rd International Conference on Machine Learning, PMLR 48:2722-2731, 2016.