

Newton-ADMM: A Distributed GPU-Accelerated Optimizer for Non-Convex Problems

K Ghosh

School of Electrical Engineering and Computer Science (EECS)

University of Ottawa

Ottawa, Canada

akademik.gk@gmail.com

October 6, 2021

1 Introduction

The critical component of a large amount of machine learning (ML) applications is estimating parameters of a model from a given data set. The process of parameter estimation problem usually convert into one of finding a minima of an appropriately suitably formulated objective or cost function of the machine learning applications. The main challenges in today's machine learning applications dealing with real big data are: a) very large numbers of parameters of the models - which is essentially equivalent to high dimensional optimization problems, b) huge training data set c) low generalization errors in learning models [1]. To address these challenges quite a size amount of efforts in research has been put. The most prevalent utilized optimization technique in machine learning problem is gradient decent and its variants, such as stochastic gradient descent (SGD). Generally gradient algorithms, that only depends on gradient information are referred to as first-order methods. The curvature information in the form of Hessian or its approximations seem to have gained improvements in performance as shown in their execution time, rate of convergence and predictions of the model in the recent research [5,6,25]. The other key challenges in machine learning optimization problems is distributed nature of the huge training data-set. As it is almost impractical to gather the whole training data set on single node or machine. Also, huge training data set cannot be processed serially on a single node due to the lack of resources, privacy of data as data can be imported to or shared with a centralized node, lessening the time for optimization. To answer these key challenges, there is a requirement for the optimization methods that are appropriately tailored to parallel and distributed computing environments. To solve the aforementioned issue, Chih-Hao Fang et al - in a recent paper - presented a new distributed optimizer for classification problems, which associates a GPU-accelerated Newton-type solver with the global consensus formulation of Alternating Direction of Method Multipliers (ADMM). For this project, I have tried to extend their work to non-convex problems generated by deep neural networks by incorporating serial non-convex solvers Newton-MR into their distributed framework.

2 Literature Review

The main reasons why first-order methods such as gradient descent and its variants are usually used in ML applications are the low cost computing cost for per-iterations and the ease of implementation [27, 25]. However, there are drawbacks too. These methods require quite a large number of iterations in order to attain generalization; mainly due to the fact that they are reactive to ill-conditioning problem i.e. it becomes hard to attain generalization. On the other hand, the second-order methods utilize curvature information by means of Hessian matrix - consequently they are not prone to ill-conditioning problem and are not affected by hyper-parameter tuning [23, 6, 28]. However, in order to process the Hessian matrix, they require higher memory and other computation resources. One way to avoid this issue is to use quasi-Newton methods to approximate the Hessian matrix using the history of gradients. But, for approximating the Hessian matrix, the gradients need to be stored and to satisfy the strong Wolfe condition, extra computation cost is added up. Furthermore, these methods are reported to be unstable when used on conjunction with mini-batches [23, 26]. Of late various distributed solvers have been proposed both first-order and second order methods. Given all of these proposed solvers for first order methods have minimum communication overhead, they have higher communication costs for large number of exchanging of messages for each mini-batch and total number of iteration [14, 20, 28, 29]. Second-order methods are proposed to reduce communication cost as well as to improve the convergence rate [18, 3, 22, 30]. To approximate Newton direction, DANE i.e. DynaNewton - Accelerating Newton's [3] and the accelerated inexact of DANE, called as AIDE [22] use Stochastic Variance Reduced Gradient (SVRG)[30] as the subproblem solver. These methods are often affected to the fine-tuning of SVRG. Another scheme name DiSCO employs distributed Preconditioned Conjugate Gradient (PCG) to approximate the Newton direction. In this scheme the total number of communications among nodes per PCG call is commensurate to the number of PCG iterations [7]. There is another proposed method named GIANT, which performs conjugate gradient (CG) on each node and approximates the Newton Direction by means of averaging the solution from each CG call [6]. The experimental results have demonstrated that GIANT performs better than DANE, AIDE and DiSCO. An adaptive approach, which is similar to trust-region methods and adapts dynamically the auxiliary model to compensate for modeling errors proposed by Dunner et al [8] is claimed to outperform GIANT but it does not do well on sparse data sets. Another recently proposed variant name DINGO [18] is derived by optimization of the gradient's norm as a surrogate function. It does not apply any specific form to the underlying functions and it can be applied beyond convexity problem - class of non-convex functions such as invex, which treats convexity as a special sub-class. In addition, it supports the arbitrary distribution of the data across the computing environment. However, it can converge to unwanted stationary points if invexity is not present. Alternating Direction Method of Multipliers (ADMM)[2] is a famous choice for distributed environment. It is based on an augmented Lagrangian framework - which tries to solve the global consensus problem. It solves such problem by alternating iterations on primal or dual variables. As a result, it inherits the advantages of the superior convergence properties of the method multipliers and the decomposability of dual ascent. It only needs one round of communication per iteration. However, the choice of the penalty parameter [31] and the local subproblem solvers affect its performance a lot.

References

- [1] Fang, C.H., Kylasa, S.B., Roosta, F., Mahoney, M.W. and Grama, A., 2020, November. Newton-ADMM: A distributed GPU-accelerated optimizer for multiclass classification problems. In SC20: International Conference for High Performance Computing, Networking, Storage and Analysis (pp. 1-12). IEEE.
- [2] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. "Distributed optimization and statistical learning via the alternating direction method of multipliers". In: Foundations and Trends® in Machine learning 3.1 (2011), pp. 1-122.
- [3] Hadi Daneshmand, Aurelien Lucchi, and Thomas Hofmann. "DynaNewton-Accelerating Newton's Method for Machine Learning". In: arXiv preprint arXiv:1605.06561 (2016).
- [4] Fred Roosta, Yang Liu, Peng Xu, and Michael W Mahoney. "Newton-MR: Newton's Method Without Smoothness or Convexity". In: arXiv preprint arXiv:1810.00303 (2018).
- [5] Sudhir B Kylasa. "HIGHER ORDER OPTIMIZATION TECHNIQUES FOR MACHINE LEARNING". PhD thesis. Purdue University Graduate School, 2019.
- [6] Farbod Roosta-Khorasani and Michael W Mahoney. Sub-sampled Newton methods. Mathematical Programming, 174(1-2):293326, 2019.
- [7] Yuchen Zhang and Xiao Lin. "DiSCO: Distributed optimization for self-concordant empirical loss". In: International conference on machine learning. 2015, pp. 362-370.
- [8] Celestine D'unner, Aurelien Lucchi, Matilde Gargiani, An Bian, Thomas Hofmann, and Martin Jaggi. "A Distributed Second-Order Algorithm You Can Trust". In: arXiv preprint arXiv:1806.07569 (2018).
- [9] L'eon Bottou, Frank E Curtis, and Jorge Nocedal. "Optimization methods for large-scale machine learning". In: arXiv preprint arXiv:1606.04838 (2016).
- [10] Hadi Daneshmand, Aurelien Lucchi, and Thomas Hofmann. "DynaNewton-Accelerating Newton's Method for Machine Learning". In: arXiv preprint arXiv:1605.06561 (2016).
- [11] Fred Roosta, Yang Liu, Peng Xu, and Michael W Mahoney. "Newton-MR: Newton's Method Without Smoothness or Convexity". In: arXiv preprint arXiv:1810.00303 (2018).?
- [12] L. Angelani, R. Di Leonardo, G. Ruocco, A. Scala, and F. Sciortino. Saddles in the Energy Landscape Probed by Supercooled Liquids. Physical review letters, 85(25):5356, 2000.
- [13] Y. Arjevani, O. Shamir, and R. Shi. Oracle Complexity of Second-Order Methods for Smooth Convex Optimization. Mathematical Programming, pages 134, 2017.?
- [14] Jianmin Chen, Xinghao Pan, Rajat Monga, Samy Bengio, and Rafal Jozefowicz. "Revisiting distributed synchronous SGD". In: arXiv preprint arXiv:1604.00981 (2016).?

- [15] Y. Bengio et al. Learning deep architectures for AI. Foundations and trends R in Machine Learning, 2(1):1127, 2009.?
- [16] S. Bellavia, C. Cartis, N. I. M. Gould, B. Morini, and Ph. L. Toint. Convergence of a regularized Euclidean residual algorithm for nonlinear least-squares. SIAM Journal on Numerical Analysis, 48(1):129, 2010.?
- [17] Groueix, T., Fisher, M., Kim, V., Russell, B., Aubry, M.: Atlasnet: A papier-m^ache approach to learning 3d surface generation. In : *CVPR2018*(2018)
- [18] Rixon Crane and Fred Roosta. "DINGO: Distributed Newton-Type Method for Gradient-Norm Optimization". In: Proceedings of the Advances in Neural Information Processing Systems. Accepted. 2019.?
- [19] D. Calvetti, B. Lewis, and L. Reichel. L-curve for the MINRES method. In Advanced Signal Processing Algorithms, Architectures, and Implementations X, volume 4116, pages 385396. International Society for Optics and Photonics, 2000.?
- [20] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. "Accurate, large minibatch SGD: training imagenet in 1 hour". In: arXiv preprint arXiv:1706.02677 (2017).?
- [21] K. v. d. Doel and U. Ascher. Adaptive and stochastic algorithms for EIT and DC resistivity problems with piecewise constant solutions and many measurements. SIAM J. Scient. Comput., 34:DOI: 10.1137/110826692 2012?
- [22] Sashank J Reddi, Jakub Konečný, Peter Richtárik, Barnabás Póczos, and Alex Smola. "AIDE: Fast and communication efficient distributed optimization". In: arXiv preprint arXiv:1608.06879 (2016).?
- [23] B. Kylasa, F. Roosta-Khorasani, M. W. Mahoney, and A. Grama. GPU Accelerated Sub-Sampled Newton's Method. arXiv preprint arXiv:1802.09113. Accepted for publication in the Proceedings of SIAM SDM 2019.?
- [24] H. Karimi, J. Nutini, and M. Schmidt. Linear convergence of gradient and proximal-gradient methods under the Polyak-Lojasiewicz condition. Joint European Conference on Machine Learning and Knowledge Discovery in Databases, pages 795811, 2016.?
- [25] Sébastien Bubeck et al. "Convex optimization: Algorithms and complexity". In: Foundations and Trends® in Machine Learning 8.3-4 (2015), pp. 231-357
- [26] Albert S Berahas, Raghu Bollapragada, and Jorge Nocedal. An Investigation of Newton-Sketch and Subsampled Newton Methods. arXiv preprint arXiv:1705.06211, 2017.
- [27] Amir Beck. First-Order Methods in Optimization. Vol. 25. SIAM, 2017
- [28] Jeffrey Dean et al. "Large scale distributed deep networks". In: Advances in neural information processing systems. 2012, pp. 1223–1231.
- [29] Peter H Jin, Qiaochu Yuan, Forrest Iandola, and Kurt Keutzer. "How to scale distributed deep learning?" In: arXiv preprint arXiv:1611.04581 (2016).

- [30] Shusen Wang, Farbod Roosta-Khorasani, Peng Xu, and Michael W Mahoney. “GIANT: Globally Improved Approximate Newton Method for Distributed Optimization”. In: Advances in Neural Information Processing Systems (NIPS). 2018, pp. 2338–2348.
- [31] Zheng Xu, Gavin Taylor, Hao Li, Mario Figueiredo, Xiaoming Yuan, and Tom Goldstein. “Adaptive consensus ADMM for distributed optimization”. In: arXiv preprint arXiv:1706.02869 (2017).