Newton Methods for Neural Networks: Introduction

Chih-Jen Lin National Taiwan University

Last updated: May 25, 2020

< □ > < □ > < □ > < □ > < □ > < □ >







Hessian and Gaussian-Newton Matrices



Chih-Jen Lin (National Taiwan Univ.)

Outline



2 Newton method

3 Hessian and Gaussian-Newton Matrices



Chih-Jen Lin (National Taiwan Univ.)

Optimization Methods Other than Stochastic Gradient

- We have explained why stochastic gradient is popular for deep learning
- The same reasons may explain why other methods are not suitable for deep learning
- But we also notice that from the simplest SG to what people are using many modifications were made
- Can we extend other optimization methods to be suitable for deep learning?



(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

Outline

Introduction



3 Hessian and Gaussian-Newton Matrices



Chih-Jen Lin (National Taiwan Univ.)

Newton Method

• Consider an optimization problem

 $\min_{\theta} f(\theta)$

 Newton method solves the 2nd-order approximation to get a direction d

$$\min_{\boldsymbol{d}} \quad \nabla f(\boldsymbol{\theta})^{\mathsf{T}} \boldsymbol{d} + \frac{1}{2} \boldsymbol{d}^{\mathsf{T}} \nabla^2 f(\boldsymbol{\theta}) \boldsymbol{d} \qquad (1)$$

• If $f(\theta)$ isn't strictly convex, (1) may not have a unique solution

Newton Method (Cont'd)

- Then (1) can be solved by

$$G\boldsymbol{d} = -\nabla f(\boldsymbol{\theta})$$

• The resulting direction is a descent one

$$abla f(oldsymbol{ heta})^T oldsymbol{d} = -
abla f(oldsymbol{ heta})^T G^{-1}
abla f(oldsymbol{ heta}) < 0$$

イロト 不得 トイヨト イヨト

Newton Method (Cont'd)

The procedure:

while stopping condition not satisfied **do** Let G be $\nabla^2 f(\theta)$ or its approximation Exactly or approximately solve

$$G \boldsymbol{d} = -\nabla f(\boldsymbol{\theta})$$

Find a suitable step size α Update

 $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \boldsymbol{d}.$

end while

< □ > < □ > < □ > < □ > < □ > < □ >

Step Size I

- Selection of the step size $\alpha :$ usually two types of approaches
 - Line search
 - Trust region (or its predecessor: Levenberg-Marquardt algorithm)
- If using line search, details are similar to what we had for gradient descent
- $\bullet\,$ We gradually reduce α such that

$$f(\boldsymbol{\theta} + \alpha \boldsymbol{d}) < f(\boldsymbol{\theta}) + \nu \nabla f(\boldsymbol{\theta})^{T} (\alpha \boldsymbol{d})$$



< ロト < 同ト < ヨト < ヨト

Newton versus Gradient Descent I

- We know they use second-order and first-order information respectively
- What are their special properties?
- It is known that using higher order information leads to faster final local convergence

・ 何 ト ・ ヨ ト ・ ヨ ト

Newton versus Gradient Descent II

• An illustration (modified from Tsai et al. (2014)) presented earlier



Slow final convergence Fast final convergence



< ロ > < 同 > < 回 > < 回 >

Newton versus Gradient Descent III

- But the question is for machine learning why we need fast local convergence?
- The answer is no
- However, higher-order methods tend to be more robust
- Their behavior may be more consistent across easy and difficult problems
- It's known that stochastic gradient is sometimes sensitive to parameters
- Thus what we hope to try here is if we can have a more robust optimization method

Difficulties of Newton for NN I

• The Newton linear system

$$Gd = -\nabla f(\theta)$$
 (2)

can be large.

 $G \in R^{n \times n}$,

where n is the total number of variables

• Thus G is often too large to be stored

Difficulties of Newton for NN II

• Evan if we can store G, calculating

$$\boldsymbol{d} = -G^{-1} \nabla f(\boldsymbol{\theta})$$

is usually very expensive

• Thus a direct use of Newton for deep learning is hopeless

< ロト < 同ト < ヨト < ヨト

Existing Works Trying to Make Newton Practical I

- Many works tried to address this issue
- Their approaches significantly vary
- I roughly categorize them to two groups
 - Hessian-free (Martens, 2010; Martens and Sutskever, 2012; Wang et al., 2020; Henriques et al., 2018)
 - Hessian approximation (Martens and Grosse, 2015; Botev et al., 2017; Zhang et al., 2017) In particular, diagonal approximation

< □ > < □ > < □ > < □ > < □ > < □ >

Existing Works Trying to Make Newton Practical II

- There are many others where I didn't put into the above two groups for various reasons (Osawa et al., 2019; Wang et al., 2018; Chen et al., 2019; Wilamowski et al., 2007)
- There are also comparisons (Chen and Hsieh, 2018)
- With the many possibilities it is difficult to reach conclusions
- We decide to first check the robustness of standard Newton methods on small-scale data



< ロト < 同ト < ヨト < ヨ

Existing Works Trying to Make Newton Practical III

• Thus in our discussion we try not to do approximations

イロト イヨト イヨト イヨト

17 / 30

Outline

Introduction

2 Newton method

Hessian and Gaussian-Newton Matrices



Chih-Jen Lin (National Taiwan Univ.)

Introduction

- We will check techniques to address the difficulty of storing or inverting the Hessian
- But before that let's derive the mathematical form

Hessian Matrix I

• For CNN, the gradient of $f(\theta)$ is

$$\nabla f(\boldsymbol{\theta}) = \frac{1}{C}\boldsymbol{\theta} + \frac{1}{I}\sum_{i=1}^{I} (J^{i})^{T} \nabla_{\boldsymbol{z}^{L+1,i}} \xi(\boldsymbol{z}^{L+1,i}; \boldsymbol{y}^{i}, \boldsymbol{Z}^{1,i}),$$
(3)

where

$$J^{i} = \begin{bmatrix} \frac{\partial z_{1}^{L+1,i}}{\partial \theta_{1}} & \cdots & \frac{\partial z_{1}^{L+1,i}}{\partial \theta_{n}} \\ \vdots & \vdots & \vdots \\ \frac{\partial z_{nL+1}^{L+1,i}}{\partial \theta_{1}} & \cdots & \frac{\partial z_{nL+1}^{L+1,i}}{\partial \theta_{n}} \end{bmatrix}_{n_{L+1} \times n}, \quad i = 1, \dots, I, \quad (4)$$

æ

イロト イヨト イヨト イヨト

Hessian Matrix II

is the Jacobian of $z^{L+1,i}(\theta)$.

• The Hessian matrix of $f(\theta)$ is

$$\nabla^{2} f(\boldsymbol{\theta}) = \frac{1}{C} \mathcal{I} + \frac{1}{I} \sum_{i=1}^{I} (J^{i})^{T} B^{i} J^{i} + \frac{1}{I} \sum_{i=1}^{I} \sum_{j=1}^{n_{L}} \frac{\partial \xi(\boldsymbol{z}^{L+1,i}; \boldsymbol{y}^{i}, \boldsymbol{Z}^{1,i})}{\partial z_{j}^{L+1,i}} \begin{bmatrix} \frac{\partial^{2} z_{j}^{L+1,i}}{\partial \theta_{1} \partial \theta_{1}} & \cdots & \frac{\partial^{2} z_{j}^{L+1,i}}{\partial \theta_{1} \partial \theta_{n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial^{2} z_{j}^{L+1,i}}{\partial \theta_{n} \partial \theta_{1}} & \cdots & \frac{\partial^{2} z_{j}^{L+1,i}}{\partial \theta_{n} \partial \theta_{n}} \end{bmatrix},$$

э

イロト イヨト イヨト イヨト

Hessian Matrix III

where \mathcal{I} is the identity matrix and B^i is the Hessian of $\xi(\cdot)$ with respect to $z^{L+1,i}$:

$$B^{i} = \nabla^{2}_{\mathbf{z}^{L+1,i},\mathbf{z}^{L+1,i}} \xi(\mathbf{z}^{L+1,i}; \mathbf{y}^{i}, Z^{1,i})$$

• More precisely,

$$B_{ts}^{i} = \frac{\partial^{2} \xi(\boldsymbol{z}^{L+1,i}; \boldsymbol{y}^{i}, \boldsymbol{Z}^{1,i})}{\partial z_{t}^{L+1,i} \partial z_{s}^{L+1,i}}, \forall t, s = 1, \dots, n_{L+1}.$$
 (5)

• Usually B^i is very simple.

< 日 > < 同 > < 回 > < 回 > .

Hessian Matrix IV

• For example, if the squared loss is used,

$$\xi(\mathbf{z}^{L+1,i};\mathbf{y}^{i}) = ||\mathbf{z}^{L+1,i} - \mathbf{y}^{i}||^{2}$$

then

$$B^{i} = \begin{bmatrix} 2 & & \\ & \ddots & \\ & & 2 \end{bmatrix}$$

• Usually we consider a convex loss function

$$\xi(\mathbf{z}^{L+1,i};\mathbf{y}^i)$$

with respect to $z^{L+1,i}$



Chih-Jen Lin (National Taiwan Univ.)

(日)

Hessian Matrix V

- Thus Bⁱ is positive semi-definite
- The last term of ∇²f(θ) may not be positive semi-definite
- Note that for a twice differentiable function $f(\theta)$ $f(\theta)$ is convex

if and only if

 $abla^2 f(oldsymbol{ heta})$ is positive semi-definite



イロト イポト イヨト イヨト

Jacobian Matrix

• The Jacobian matrix of $m{z}^{L+1,i}(m{ heta})\in R^{n_{L+1}}$ is

$$J^{i} = \begin{bmatrix} \frac{\partial z_{1}^{L+1,i}}{\partial \theta_{1}} & \cdots & \frac{\partial z_{1}^{L+1,i}}{\partial \theta_{n}} \\ \vdots & \vdots & \vdots \\ \frac{\partial z_{n_{L}}^{L+1,i}}{\partial \theta_{1}} & \cdots & \frac{\partial z_{n_{L}}^{L+1,i}}{\partial \theta_{n}} \end{bmatrix} \in R^{n_{L+1} \times n}, \ i = 1, \dots I.$$

• n_{L+1} : # of neurons in the output layer

- n: number of total variables
- $n_{L+1} \times n$ can be large

イロン イ理 とく ヨン イヨン

Gauss-Newton Matrix I

- The Hessian matrix ∇²f(θ) is now not positive definite.
- We may need a positive definite approximation
- This is a deep research issue
- Many existing Newton methods for NN has considered the Gauss-Newton matrix (Schraudolph, 2002)

$$G = \frac{1}{C}\mathcal{I} + \frac{1}{I}\sum_{i=1}^{I} (J^i)^T B^i J^i$$

by removing the last term in $\nabla^2 f(\theta)$

Gauss-Newton Matrix II

- The Gauss-Newton matrix is positive definite if *Bⁱ* is positive semi-definite
- This can be achieved if we use a convex loss function in terms of z^{L+1,i}(θ)
- We then solve

$$G\boldsymbol{d} = -\nabla f(\boldsymbol{\theta})$$

イロト イヨト イヨト イヨト

References I

- A. Botev, H. Ritter, and D. Barber. Practical Gauss-Newton optimisation for deep learning. In Proceedings of the 34th International Conference on Machine Learning, pages 557–565, 2017.
- P. H. Chen and C.-J. Hsieh. A comparison of second-order methods for deep convolutional neural networks, 2018. URL https://openreview.net/forum?id=HJYoqzbC-.
- S.-W. Chen, C.-N. Chou, and E. Y. Chang. An approximate second-order method for training fully-connected neural networks. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence*, 2019.
- J. F. Henriques, S. Ehrhardt, S. Albanie, and A. Vedaldi. Small steps and giant leaps: Minimal Newton solvers for deep learning, 2018. arXiv preprint 1805.08095.
- J. Martens. Deep learning via Hessian-free optimization. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, 2010.
- J. Martens and R. Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 2408–2417, 2015.
- J. Martens and I. Sutskever. Training deep and recurrent networks with Hessian-free optimization. In *Neural Networks: Tricks of the Trade*, pages 479–535. Springer, 2012.



< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 >

References II

- K. Osawa, Y. Tsuji, Y. Ueno, A. Naruse, R. Yokota, and S. Matsuoka. Large-scale distributed second-order optimization using kronecker-factored approximate curvature for deep convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision* and Pattern Recognition (CVPR), 2019.
- N. N. Schraudolph. Fast curvature matrix-vector products for second-order gradient descent. *Neural Computation*, 14(7):1723–1738, 2002.
- C.-H. Tsai, C.-Y. Lin, and C.-J. Lin. Incremental and decremental training for linear classification. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2014. URL http://www.csie.ntu.edu.tw/~cjlin/papers/ws/inc-dec.pdf.
- C.-C. Wang, K.-L. Tan, C.-T. Chen, Y.-H. Lin, S. S. Keerthi, D. Mahajan, S. Sundararajan, and C.-J. Lin. Distributed Newton methods for deep learning. *Neural Computation*, 30(6): 1673–1724, 2018. URL http://www.csie.ntu.edu.tw/~cjlin/papers/dnn/dsh.pdf.
- C.-C. Wang, K. L. Tan, and C.-J. Lin. Newton methods for convolutional neural networks. ACM Transactions on Intelligent Systems and Technology, 2020. URL https://www.csie.ntu.edu.tw/~cjlin/papers/cnn/newton-CNN.pdf. To appear.
- B. M. Wilamowski, N. Cotton, J. Hewlett, and O. Kaynak. Neural network trainer with second order learning algorithms. In *In Proceedings of the 11th International Conference on Intelligent Engineering Systems*, 2007.

References III

H. Zhang, C. Xiong, J. Bradbury, and R. Socher. Block-diagonal Hessian-free optimization for training neural networks, 2017. arXiv preprint arXiv:1712.07296.