Block Mean Approximation for Efficient Second Order Optimization

Yao Lu, Mehrtash Harandi, Richard Hartley, Razvan Pascanu

ANU & Data61 & Google DeepMind

November 25, 2023

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Optimization

For $f(\boldsymbol{\theta}) : \mathbb{R}^n \to \mathbb{R}$

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

◆□ ▶ < 圖 ▶ < 圖 ▶ < 圖 ▶ < 圖 • 의 Q @</p>

Optimization

For $f(\boldsymbol{\theta}) : \mathbb{R}^n \to \mathbb{R}$

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

First order method

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$$

Second order method

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \mathbf{G}^{-1} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$$

Optimization

For $f(\boldsymbol{\theta}) : \mathbb{R}^n \to \mathbb{R}$

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

First order method

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$$

Second order method

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \mathbf{G}^{-1} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$$

- Newton method
- Natural gradient
- Adaptive gradient

Example: logistic regression

Logistic regression example, with $n=500,\,p=100:$ we compare gradient descent and Newton's method, both with backtracking



k

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ⊙

Second Order Optimization

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \mathbf{G}^{-1} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Compute \mathbf{G}^{-1} is expensive.

Approximation is needed.

Approximation of $\boldsymbol{\mathsf{G}}$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \mathbf{G}^{-1} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$$



Approximation of ${\boldsymbol{\mathsf{G}}}$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \mathbf{G}^{-1} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$$

- Diagonal
- Block diagonal
- Low rank
- Kronecker product



(a) Original



(b) Approximate



▲ロト ▲圖ト ▲ヨト ▲ヨト 三ヨー のへで



・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト ・ ヨ

each block.







It can be computed by inverting a matrix of size 🛄 and 🚺

Theorem

For the non-singular matrix $\overline{\Lambda} + \overline{B}$, where $\overline{\Lambda}$ and \overline{B} are the diagonal and full expansion of Λ and B with respect to the partition vector s,

$$(\bar{\mathbf{\Lambda}} + \bar{\mathbf{B}})^{-1} = \bar{\mathbf{\Lambda}}^{-1} + \bar{\mathbf{D}}$$
(1)

where $\bar{\mathbf{D}}$ is the full expansion matrix with partition vector \mathbf{s} of

$$\mathbf{D} = (\mathbf{\Lambda}\mathbf{S} + \mathbf{S}\mathbf{B}\mathbf{S})^{-1} - (\mathbf{\Lambda}\mathbf{S})^{-1}$$
(2)

where $\mathbf{S} = \operatorname{diag}(\mathbf{s})$.

Theorem

For the non-singular matrix $\overline{\Lambda} + \overline{B}$, where $\overline{\Lambda}$ and \overline{B} are the diagonal and full expansion of Λ and B with respect to the partition vector s,

$$(\bar{\mathbf{\Lambda}} + \bar{\mathbf{B}})^{-\frac{1}{2}} = \bar{\mathbf{\Lambda}}^{-\frac{1}{2}} + \bar{\mathbf{D}}$$
(3)

where $\bar{\mathbf{D}}$ is the full expansion matrix with partition vector \mathbf{s} of

$$\mathbf{D} = \mathbf{S}^{-\frac{1}{2}} \left[(\mathbf{\Lambda} + \mathbf{S}^{\frac{1}{2}} \mathbf{B} \mathbf{S}^{\frac{1}{2}})^{-\frac{1}{2}} - \mathbf{\Lambda}^{-\frac{1}{2}} \right] \mathbf{S}^{-\frac{1}{2}}$$
(4)

where $\mathbf{S} = diag(\mathbf{s})$.



(c) Original



(d) Approximate

Optimal Block Mean Approximation

Proposition

The optimal block mean approximation of ${\bf M}$ with the partition vector ${\bf s}$ according to the Frobenius norm

$$\min_{\bar{\mathbf{\Lambda}},\bar{\mathbf{B}}} \|\bar{\mathbf{\Lambda}} + \bar{\mathbf{B}} - \mathbf{M}\|_{F}^{2}$$
(5)

is given by

$$b_{ij} = \begin{cases} 0, & i = j, s_i = 1, \\ \frac{\sum_{mn} M_{mn}^{ii} - \sum_m M_{mm}^{ii}}{s_i(s_i - 1)}, & i = j, s_i \neq 1, \\ \frac{\sum_{mn} M_{mn}^{ij}}{s_i s_j}, & i \neq j, \end{cases}$$
(6)
$$\lambda_i = \frac{1}{s_i} \sum_m M_{mm}^{ii} - b_{ii} .$$
(7)



For neural nets, each block can represent the weights in a layer.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

AdaGrad

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathbf{H}_t^{-1/2} \mathbf{g}_t$$

$$\widehat{\mathbf{H}}_t = \mathbf{Z}_t \mathbf{F}_t \mathbf{Z}_t \approx \mathbf{H}_t$$

where Z_t is diagonal and F_t is a block mean approximation matrix.

Experiments

Table 1: Small model

Conv 3x3, 3 Max Pooling 2x2 Fully Connected, 10 Softmax, 10

Table 2: Large model

Conv 3x3, 32 Conv 3x3, 32 Conv 3x3. 32 Conv 3x3, 32 Max Pooling 2x2 Conv 3x3, 32 Conv 3x3, 32 Conv 3x3. 32 Conv 3x3, 32 Max Pooling 2x2 Conv 3x3. 32 Conv 3x3, 32 Conv 3x3. 32 Conv 3x3, 32 Max Pooling 2x2 Conv 3x3, 32 Conv 3x3, 32 Conv 3x3, 32 Conv 3x3. 32 Max Pooling 2x2 Fully Connected, 10 Softmax, 10 = nac

Experiments



(g) CIFAR-10, small model (h) CIFAR-10, small model



æ

Open Questions

The right block structure?



 P_1WP_2 in Analytic Study of Families of Spurious Minima in
Two-Layer ReLU Neural Networks,NeurIPS 2021

Other applications (e.g. Gaussian processes)?