

Optimizing Neural Networks with Kronecker-Factored Approximate Curvature

Josh Romoff & Riashat Islam
IFT6085 Paper Presentation

Reasoning and Learning Lab



28th February 2018

Training Neural Networks

While deep networks have rich modelling capacity, the complex dependency between parameters can make learning difficult

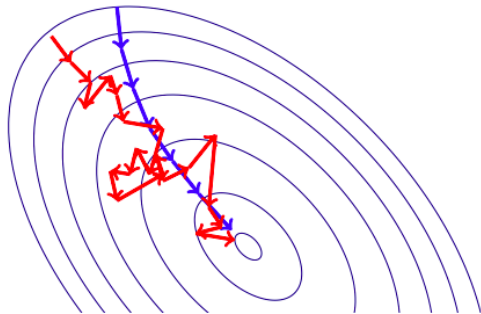
Training Neural Networks

This calls for better optimization techniques for training deep nets

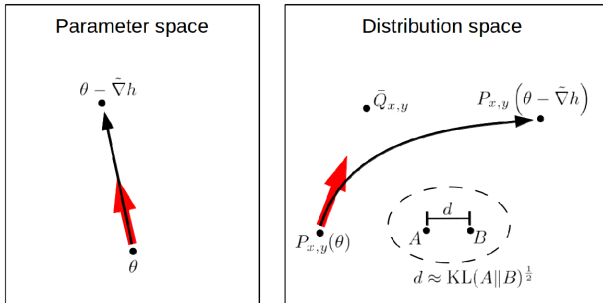
Overview

- ▶ Natural Gradients as alternative to SGD
- ▶ Fisher Information Matrix
- ▶ Approximations - KFAC
- ▶ Experimental Results
- ▶ Discussion

Stochastic Gradient Descent



Natural Gradient Descent



The Natural Gradient

$$\begin{aligned} & \arg \min_{\Delta\theta} \mathcal{L}(\theta + \Delta\theta) \\ \text{s. t. } & \mathbb{E}_{\mathbf{x} \sim \tilde{q}(\mathbf{x})} [KL(p_{\theta}(\mathbf{t}|\mathbf{x}) || p_{\theta+\Delta\theta}(\mathbf{t}|\mathbf{x}))] = \text{const.} \end{aligned}$$

- ▶ We use the expected value of the KL divergence between $p_{\theta}(t|x)$ and $p_{\theta+\delta\theta}(t|x)$ as a measure of the functional behaviour of $p_{\theta}(t|x)$ for different values of x

The Natural Gradient

Intuitively...

- ▶ Relies on the KL divergence between iterates
- ▶ Want to guarantee that the distribution of the new network will be similar to the old one
- ▶ Add a KL divergence constraint to the update, to ensure the new network will behave relatively similar to the old one

Why Natural Gradients

The KL divergence constraint ensure that we move along the functional manifold with constant speed, without being slowed by the curvature. It is a measure of how the probability density function changes, regardless of how it is parameterized.

- ▶ In mini-batch updates, some outliers may make drastic changes to network's parameters
- ▶ Stable learning process
- ▶ Natural gradients ensure monotonic parameter improvements

*Learning locally robust to re-parameterizations of the model
→ the functional behaviour of p does not depend on the parameterization of the model*

Why Natural Gradients

The gradients are of p_θ which acts as a proxy for the cost L . The KL constraint on the probability distributions ensures

- ▶ Each gradient update will make ϵ change to the model
- ▶ The model does not change by more than ϵ

This can provide some kind of robustness to overfitting.

- ▶ The model is not allowed to move too far in some direction, if moving along that direction changes the density computed by the model substantially.

Fisher Information Matrix

The metric for natural gradients is determined by the Fisher Information Matrix

$$F = \mathbb{E} \left[\frac{d \log p(y|x, \theta)}{d\theta} \frac{d \log p(y|x, \theta)}{d\theta}^\top \right] = \mathbb{E}[\mathcal{D}\theta \mathcal{D}\theta^\top]$$

ie, the covariance of the gradients of the model log probabilities w.r.t its parameters

The natural gradient is the direction obtained by

$$\nabla_N = F_\theta^{-1} \nabla \quad (1)$$

Difficulties with Natural Gradient

$$\nabla_N L(\theta) = \nabla L(\theta) F^{-1} \quad (2)$$

where F is the inverse of the Fisher information matrix.

- ▶ For large networks with millions of parameters, computing the inverse F^{-1} is computationally impractical.

We resort to approximation methods for efficiently computing the Fisher matrix F^{-1} and the natural gradient

- ▶ A key ingredient to designing optimization algorithms based on the natural gradients

Here, we will introduce the **KFAC** approximation method for efficiently computing the inverse of Fisher Information Matrix.

Kronecker Factored Fisher Approximation

We first note that the Fisher Matrix can be viewed as an ℓ by ℓ block matrix, where ℓ is the number of layers.

$$\begin{aligned} F &= \mathbb{E} [\mathcal{D}\theta\mathcal{D}\theta^\top] \\ &= \mathbb{E} \left[\left[\text{vec}(\mathcal{D}W_1)^\top \text{vec}(\mathcal{D}W_2)^\top \cdots \text{vec}(\mathcal{D}W_\ell)^\top \right]^\top \left[\text{vec}(\mathcal{D}W_1)^\top \text{vec}(\mathcal{D}W_2)^\top \cdots \text{vec}(\mathcal{D}W_\ell)^\top \right] \right] \\ &= \begin{bmatrix} \mathbb{E} \left[\text{vec}(\mathcal{D}W_1) \text{vec}(\mathcal{D}W_1)^\top \right] & \mathbb{E} \left[\text{vec}(\mathcal{D}W_1) \text{vec}(\mathcal{D}W_2)^\top \right] & \cdots & \mathbb{E} \left[\text{vec}(\mathcal{D}W_1) \text{vec}(\mathcal{D}W_\ell)^\top \right] \\ \mathbb{E} \left[\text{vec}(\mathcal{D}W_2) \text{vec}(\mathcal{D}W_1)^\top \right] & \mathbb{E} \left[\text{vec}(\mathcal{D}W_2) \text{vec}(\mathcal{D}W_2)^\top \right] & \cdots & \mathbb{E} \left[\text{vec}(\mathcal{D}W_2) \text{vec}(\mathcal{D}W_\ell)^\top \right] \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{E} \left[\text{vec}(\mathcal{D}W_\ell) \text{vec}(\mathcal{D}W_1)^\top \right] & \mathbb{E} \left[\text{vec}(\mathcal{D}W_\ell) \text{vec}(\mathcal{D}W_2)^\top \right] & \cdots & \mathbb{E} \left[\text{vec}(\mathcal{D}W_\ell) \text{vec}(\mathcal{D}W_\ell)^\top \right] \end{bmatrix} \end{aligned}$$

Kronecker Factored Fisher Approximation

Looking at one of the blocks, $F_{i,j}$, and by doing some re-arranging we get:

$$\begin{aligned} F_{i,j} &= \text{E} [\text{vec}(\mathcal{D}W_i) \text{vec}(\mathcal{D}W_j)^\top] = \text{E} [(\bar{a}_{i-1} \otimes g_i)(\bar{a}_{j-1} \otimes g_j)^\top] = \text{E} [(\bar{a}_{i-1} \otimes g_i)(\bar{a}_{j-1}^\top \otimes g_j^\top)] \\ &= \text{E} [\bar{a}_{i-1} \bar{a}_{j-1}^\top \otimes g_i g_j^\top] \end{aligned}$$

where we recall that $g_i = \mathcal{D}s_i$, $a_i = \phi_i(s_i)$ and $s_i = W_i a_{i-1}$.

So, each block can be viewed as the expected Kronecker product between $\bar{A}_{i,j} = \bar{a}_{i-1} \bar{a}_{j-1}^\top$ and $\bar{G}_{i,j} = \bar{g}_i \bar{g}_j^\top$

Kronecker Factored Fisher Approximation

Now for the main approximation. Consider again one of the blocks, $F_{i,j}$, we assume the following:

$$F_{i,j} = \mathbb{E} [\bar{a}_{i-1} \bar{a}_{j-1}^\top \otimes g_i g_j^\top] \approx \mathbb{E} [\bar{a}_{i-1} \bar{a}_{j-1}^\top] \otimes \mathbb{E} [g_i g_j^\top] = \bar{A}_{i-1,j-1} \otimes G_{i,j} = \tilde{F}_{i,j}$$

So, we assume that the expected Kronecker product between two matrices is equal to the Kronecker product of the expectations; i.e. that they are independent.

Further Approximations for Efficient Inversion

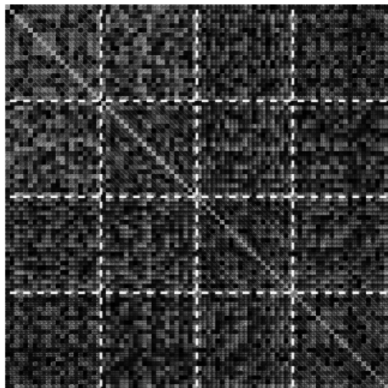
One can trade-off precision for computational savings (regarding the matrix inverse) by doing one of the following:

- ▶ Approximating \tilde{F}^{-1} as block-diagonal
- ▶ Approximating \tilde{F}^{-1} as block-tridiagonal

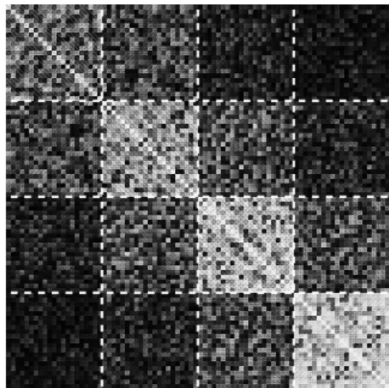
Intuitively, the block-diagonal approximation will be less precise but at substantial computational savings.

Experimental Results (Approximating \tilde{F}^{-1})

\tilde{F}

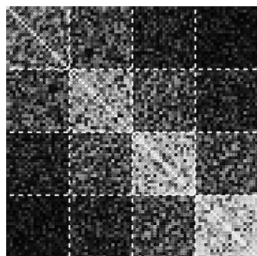
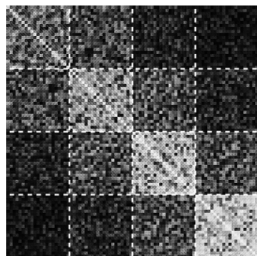


\tilde{F}^{-1}

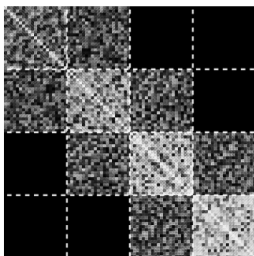
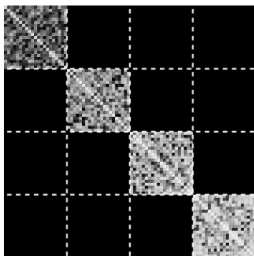


Experimental Results (Approximating \tilde{F}^{-1})

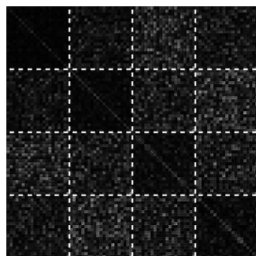
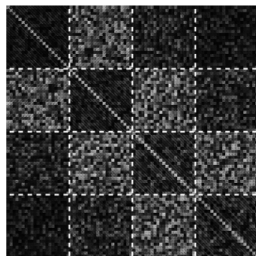
\tilde{F}^{-1}



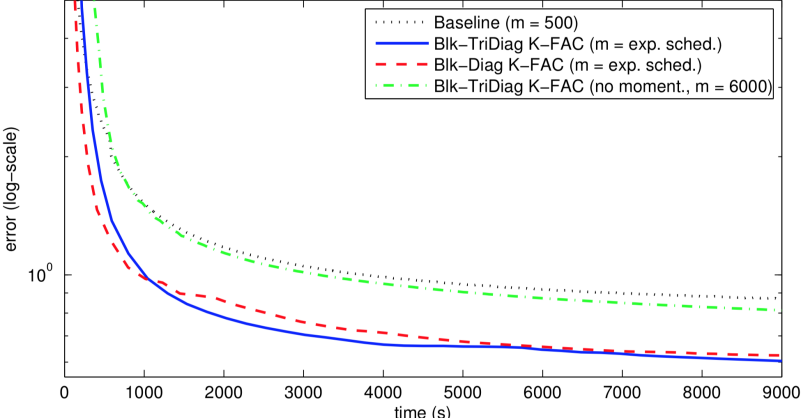
Diag,
Triag



$|\tilde{F}^{-1} - \text{Diag}|,$
 $|\tilde{F}^{-1} - \text{Triag}|$



Experimental Results (MNIST)



Summary

- ▶ Natural gradients as optimization techniques might be better
- ▶ But difficulty with the Fisher Information Matrix
- ▶ Need ways to approximate the Fisher matrix
- ▶ KFAC approximations to the rescue!

Thank You

*The only stupid question is the one you were afraid to ask
but never did*

- Rich Sutton