

Fast Approximate Natural Gradient Descent in a Kronecker-factored Eigenbasis

Thomas George^{*1}, César Laurent^{*1}, Xavier Bouthillier¹, Nicolas Ballas², Pascal Vincent^{1,2,3}

* equal contribution

Overview

Preconditioned gradient descent aims at accounting for the curvature of the optimized function. Most preconditioners require inverting a large matrix, which is unfeasible for 15M parameters neural networks.

In this work, we build upon Martens and Grosse (2015) to improve the K-FAC approximation used to efficiently invert the Fisher Information Matrix. We show that our method allows to efficiently track a more accurate approximate than K-FAC, allowing for a speedup in optimization.

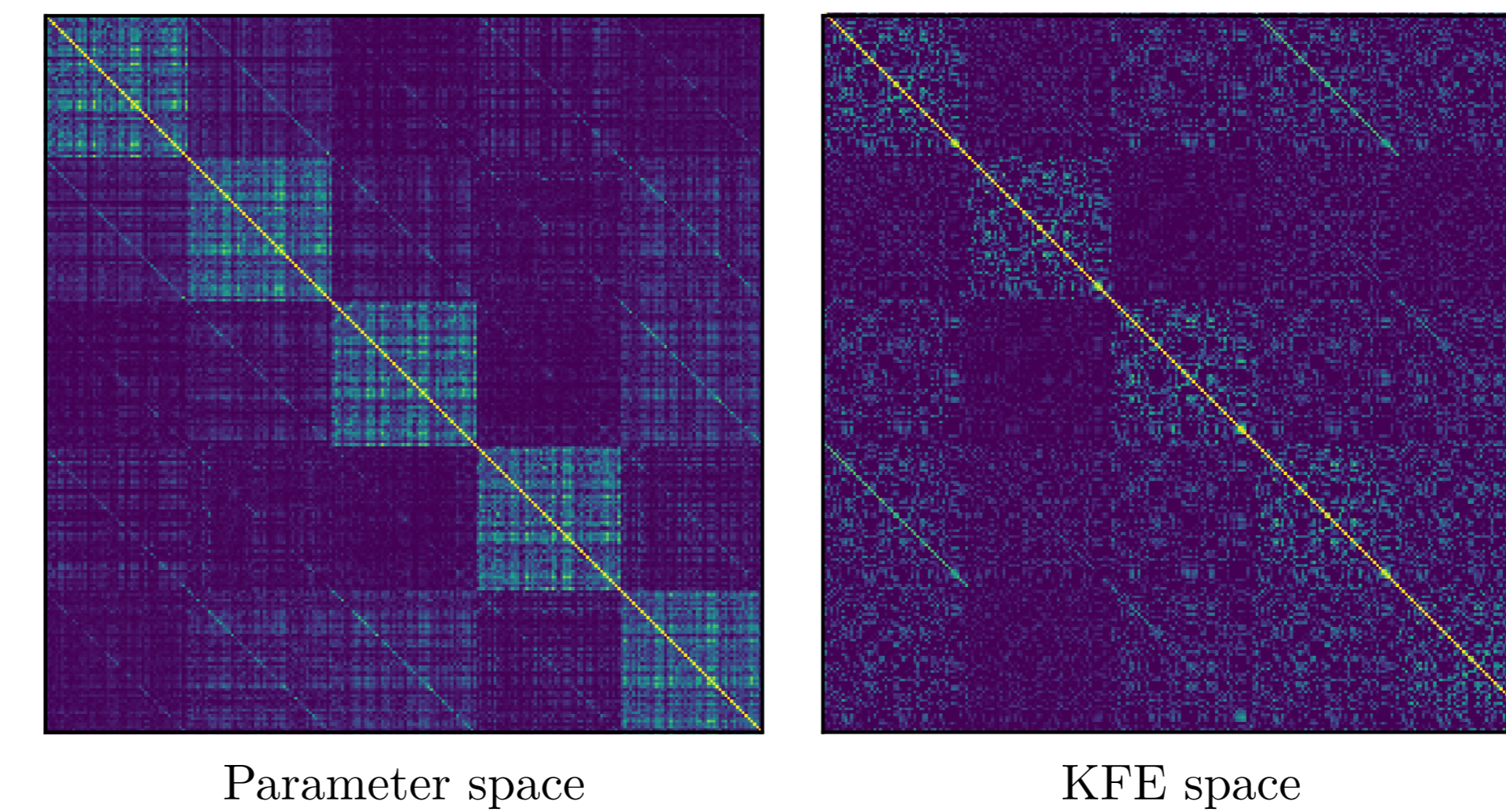
Kronecker-factored Eigenbasis (KFE)

if $A \otimes B \approx G$ then the SVD of $A \otimes B$ must be a good approximate of the SVD of G .

$$\left. \begin{aligned} A &= U_A \Lambda_A U_A^T \\ B &= U_B \Lambda_B U_B^T \end{aligned} \right\} \text{SVD of } A \text{ and } B$$

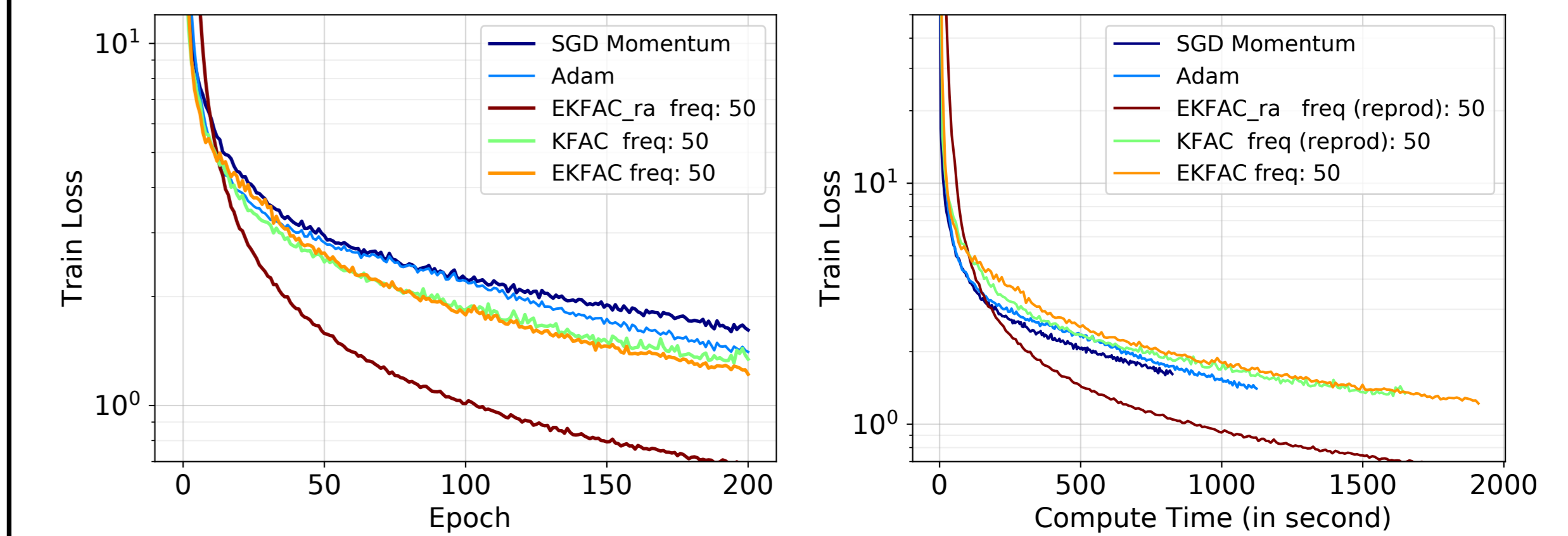
$$A \otimes B = \underbrace{(U_A \otimes U_B)}_U (\Lambda_A \otimes \Lambda_B) (U_A \otimes U_B)^T$$

approximate eigenbasis for G
= KFE



Correlation of the gradients of a single layer for a MLP on MNIST

Experiment: MNIST Deep Auto-Encoder



▷ Both EKFAC and EKFAC_ra show an optimization benefit compared to amortized K-FAC and other baselines
▷ This transfers to faster training (in wall clock time) for EKFAC_ra

Prerequisite: Preconditioned gradient descent

$$\Delta_{GD} = -\lambda \mathbb{E}[\nabla_{\theta}]$$

λ : Learning rate

$$\Delta_{PreGD} = -\lambda (G + \epsilon \mathbf{I})^{-1} \mathbb{E}[\nabla_{\theta}]$$

ϵ : Tikhonov damping

G can be:

Hessian

Second moment of gradients

Gauss-Newton

Fisher Information Matrix

"empirical Fisher"

Eigenvalue corrected K-FAC (EKFAC)

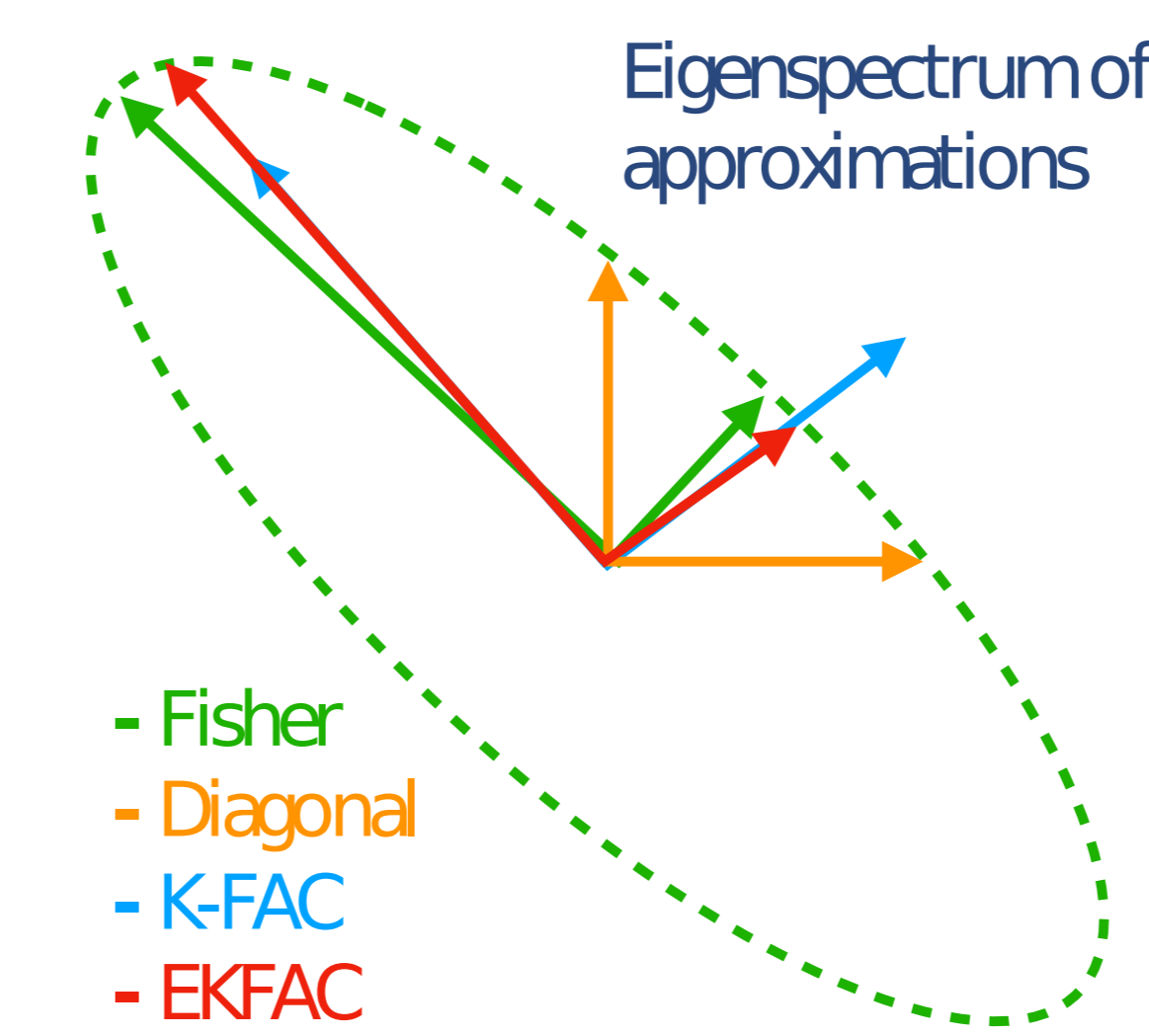
Proposal: use the KFE but rescale the eigenvalues so that they exactly match the diagonal in that basis.

$$\Lambda = \text{diag}(U^T G U)$$

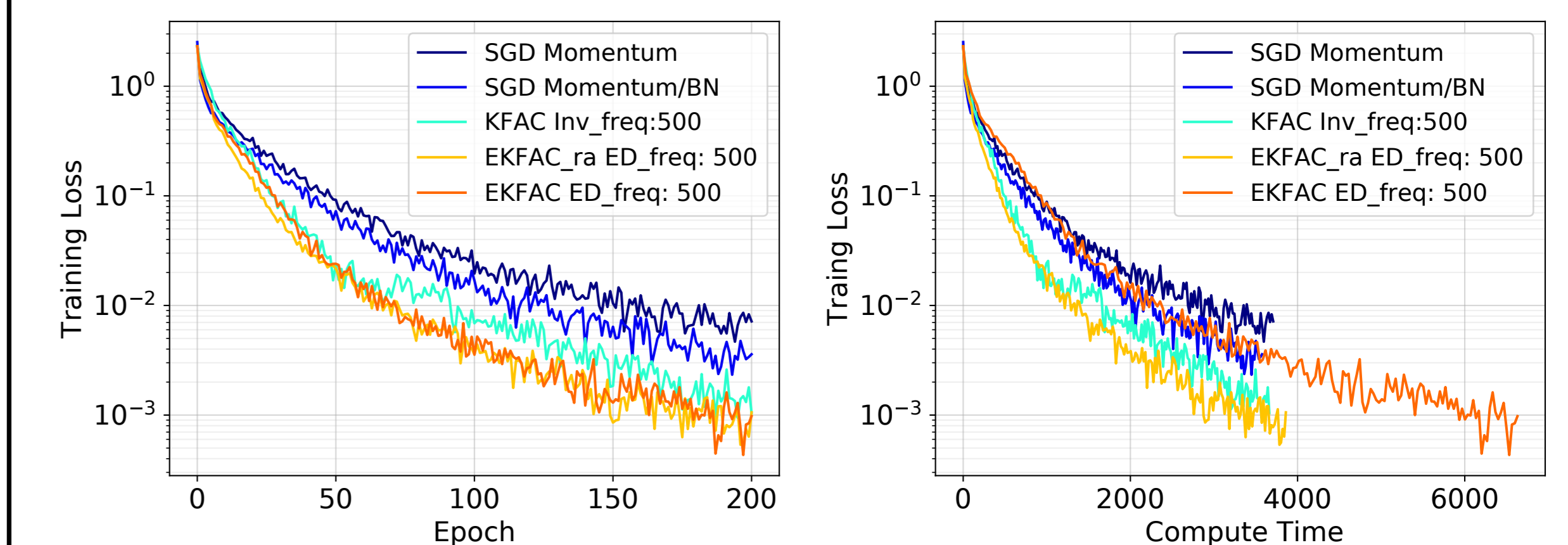
diag keeps the diagonal values and puts all other elements to 0.

Theorem: It is the optimal diagonal scaling in this basis in the sense that it minimizes $\|G - U \Lambda U^T\|_F$

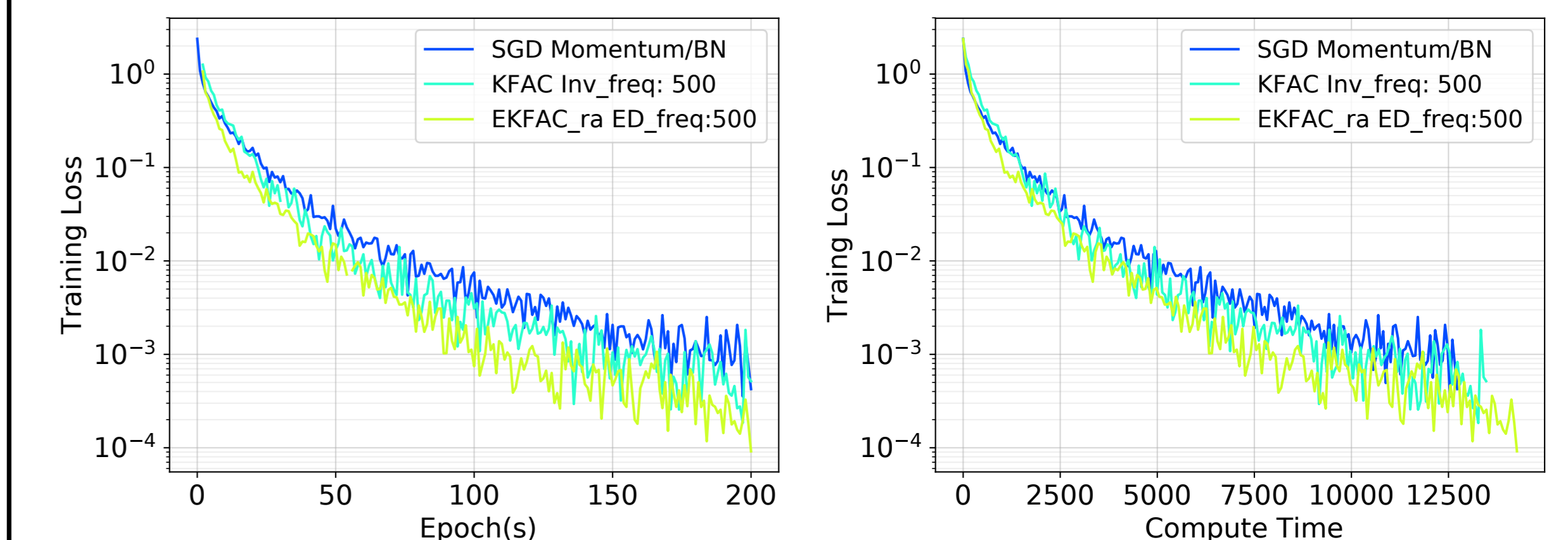
Corollary: $\|G - G_{EKFAC}\|_F \leq \|G - G_{KFAC}\|_F$



Experiment: CIFAR10 VGG11



Experiment: CIFAR10 Resnet34



▷ We did not find one set of hyperparameter for which EKFAC is below K-FAC for all epochs (and vice versa)
▷ However, if we do the model selection for each epoch, the best EKFAC will always outperform the best K-FAC
▷ K-FAC and EKFAC are very sensitive to the learning rate and the Tikhonov damping hyperparameters.

Prerequisite: Block-diagonal approximation

$$G_l = \mathbb{E}[\nabla_{\theta_l} \nabla_{\theta_l}^T] = \mathbb{E}[(x_l \otimes g_l)(x_l \otimes g_l)^T]$$

x_l : Activation from previous layer
 g_l : Gradient with respect to outgoing preactivation $\frac{\partial \ell}{\partial u_i}$

From now on we focus on a single layer and we drop the subscript l

Parameter update

$$(G + \epsilon \mathbf{I})^{-1} \mathbb{E}[\nabla_{\theta}] = U (\Lambda + \epsilon \mathbf{I})^{-1} U^T \mathbb{E}[\nabla_{\theta}]$$

$$\begin{aligned} \Lambda &= \text{diag}(U^T \mathbb{E}[\nabla_{\theta} \nabla_{\theta}^T] U) \\ &= \mathbb{E}[\text{diag}(U^T \nabla_{\theta} \nabla_{\theta}^T U)] \\ &= \mathbb{E}[(U^T \nabla_{\theta})^2] \end{aligned}$$

$\mathbb{E}[U^T \nabla_{\theta}]$
 $U^T \nabla_{\theta} = \tilde{\nabla}_{\theta}$
gradient projected in the KFE

2 interpretations:

EKFAC is a rescaled K-FAC preconditioner in the parameter space

EKFAC is a diagonal method in the KFE

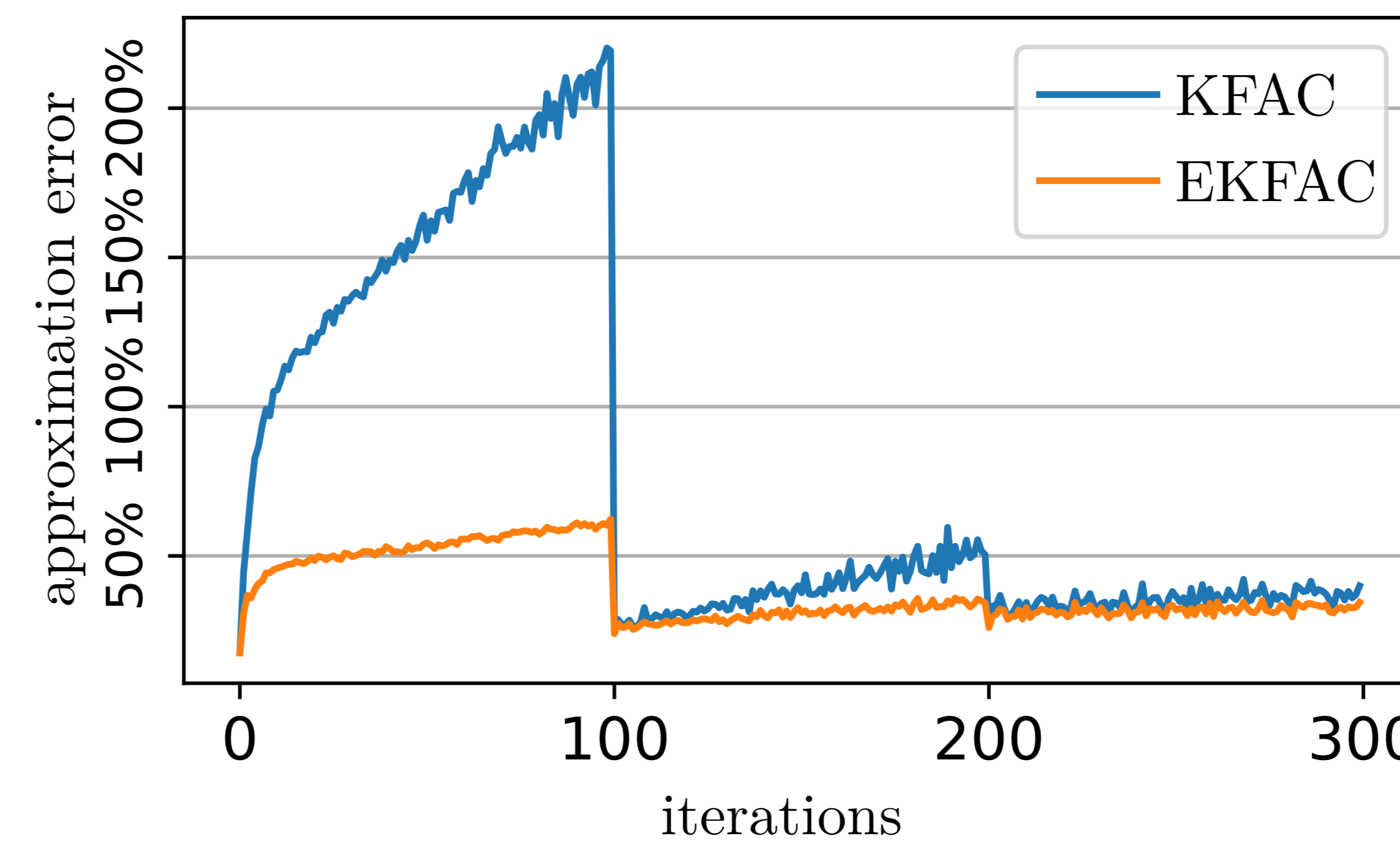
Advantages of EKFAC over K-FAC

▷ Eigendecompositions can be amortized every n updates, then Λ is cheap to compute:

- using intra minibatch 2nd moment
- using a running average estimate

▷ Allows go keep a better approximate during training compared to amortized K-FAC

Approximation error: $\frac{\|G - \hat{G}\|_F}{\|G\|_F}$



Prerequisite: K-FAC (Martens and Grosse 2015, Heskes 2000)

$$G = \mathbb{E}[xx^T \otimes gg^T] = \mathbb{E}[xx^T] \otimes \mathbb{E}[gg^T] + R$$

approximation error:

$$R = \mathbb{E}[(xx^T - \mathbb{E}[xx^T]) \otimes (gg^T - \mathbb{E}[gg^T])]$$

Kronecker product property:

$$(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$$

$mn \times mn$ $m \times m$ $n \times n$

Our PyTorch implementation of K-FAC and EKFAC is available!
<https://github.com/Thrandis/EKFAC-pytorch>

Conclusion

- ▷ EKFAC can optimize networks fast successfully without BN
- ▷ Future work: apply other diagonal methods in the KFE (RMSProp, SignSGD, ...)
- ▷ Develop regularization and improve robustness to hyperparameters

References

James Martens and Roger Grosse. "Optimizing neural networks with kronecker-factored approximate curvature." ICML 2015.
Tom Heskes. "On natural learning and pruning in multilayered perceptrons." Neural Computation, 2000