

PRECONDITIONER ON MATRIX LIE GROUP FOR SGD

Xi-Lin Li

GMEMS Technologies and Spectimbre
366 Fairview Way, Milpitas, CA 95035
lixilinx@gmail.com

ABSTRACT

We study two types of preconditioners and preconditioned stochastic gradient descent (SGD) methods in a unified framework. We call the first one the Newton type due to its close relationship to the Newton method, and the second one the Fisher type as its preconditioner is closely related to the inverse of Fisher information matrix. Both preconditioners can be derived from one framework, and efficiently estimated on any matrix Lie groups designated by the user using natural or relative gradient descent minimizing certain preconditioner estimation criteria. Many existing preconditioners and methods, e.g., RMSProp, Adam, KFAC, equilibrated SGD, batch normalization, etc., are special cases of or closely related to either the Newton type or the Fisher type ones. Experimental results on relatively large scale machine learning problems are reported for performance study.

1 INTRODUCTION

This paper investigates the use of preconditioner for accelerating gradient descent, especially in large scale machine learning problems. Stochastic gradient descent (SGD) and its variations, e.g., momentum (Rumelhart et al., 1986; Nesterov, 1983), RMSProp and Adagrad (John et al., 2011), Adam (Kingma & Ba, 2015), etc., are popular choices due to their simplicity and wide applicability. These simple methods do not use well normalized step size, could converge slow, and might involve more controlling parameters requiring fine tweaking. Convex optimization is a well studied field (Boyd & Vandenberghe, 2004). Many off-the-shelf methods there, e.g., (nonlinear) conjugate gradient descent, quasi-Newton methods, Hessian-free optimizations, etc., can be applied to small and middle scale machine learning problems without much modifications. However, these convex optimization methods may have difficulty in handling gradient noise and scaling up to problems with hundreds of millions of free parameters. For a large family of machine learning problems, natural gradient with the Fisher information metric is equivalent to a preconditioned gradient using inverse of the Fisher information matrix as the preconditioner (Amari, 1998). Natural gradient and its variations, e.g., Kronecker-factored approximate curvature (KFAC) (Martens & Grosse, 2015) and the one in (Povey et al., 2015), all use such preconditioners. Other less popular choices are the equilibrated preconditioner (Dauphin et al., 2015) and the one proposed in (Li, 2018). Momentum or the heavy-ball method provides another independent way to accelerate converge (Nesterov, 1983; Rumelhart et al., 1986). Furthermore, momentum and preconditioner can be combined to further accelerate convergence as shown in Adam (Kingma & Ba, 2015).

This paper groups the above mentioned preconditioners and preconditioned SGD methods into two classes, the Newton type and the Fisher type. The Newton type is closely related to the Newton method, and is suitable for general purpose optimizations. The Fisher type preconditioner relates to the inverse of Fisher information matrix, and is limited to a large subclass of stochastic optimization problems where the Fisher information metric can be well defined. Both preconditioners can be derived from one framework, and estimated on any matrix Lie groups designated by the user with almost the same natural or relative gradient descent methods minimizing specific preconditioner estimation criteria.

2 BACKGROUND

2.1 NOTATIONS

We consider the minimization of cost function

$$f(\boldsymbol{\theta}) = E_z[\ell(\boldsymbol{\theta}, \mathbf{z})] \quad (1)$$

where E_z takes expectation over random variable \mathbf{z} , ℓ is a loss function, and $\boldsymbol{\theta}$ is the model parameter vector to be optimized. For example, in a classification problem, ℓ could be the cross entropy loss, \mathbf{z} is a pair of input feature vector and class label, vector $\boldsymbol{\theta}$ consists of all the trainable parameters in the classification model, and E_z takes average over all samples from the training data set. By assuming second order differentiable model and loss, we could approximate $\ell(\boldsymbol{\theta}, \mathbf{z})$ as a quadratic function of $\boldsymbol{\theta}$ within a trust region around $\boldsymbol{\theta}$, i.e., $\ell(\boldsymbol{\theta}, \mathbf{z}) = \mathbf{b}_z^T \boldsymbol{\theta} + 0.5 \boldsymbol{\theta}^T \mathbf{H}_z \boldsymbol{\theta} + a_z$, where a_z is the sum of higher order approximation errors and constant term independent of $\boldsymbol{\theta}$, \mathbf{H}_z is a symmetric matrix, and subscript z in \mathbf{b}_z , \mathbf{H}_z and a_z reminds us that these three terms depend on \mathbf{z} . Clearly, these three terms depend on $\boldsymbol{\theta}$ as well, although we do not explicitly show this dependence to simplify our notations since we just consider parameter updates in the same trust region. Now, we may rewrite (1) as $f(\boldsymbol{\theta}) = \mathbf{b}^T \boldsymbol{\theta} + 0.5 \boldsymbol{\theta}^T \mathbf{H} \boldsymbol{\theta} + a$, where $\mathbf{b} = E_z[\mathbf{b}_z]$, $\mathbf{H} = E_z[\mathbf{H}_z]$, and $a = E_z[a_z]$. We do not impose any assumption, e.g., positive definiteness, on \mathbf{H} except for being symmetric. Thus, the quadratic surface in the trust region could be non-convex. To simplify our notations, we no longer consider the higher order approximation errors included in a , and simply assume that $f(\boldsymbol{\theta})$ is a quadratic function of $\boldsymbol{\theta}$ in the trust region.

2.2 PRECONDITIONER

Let us consider a certain iteration. Preconditioned SGD updates $\boldsymbol{\theta}$ as

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \mu \mathbf{P} \partial \hat{f}(\boldsymbol{\theta}) / \partial \boldsymbol{\theta} \quad (2)$$

where $\mu > 0$ is the step size, $\hat{f}(\boldsymbol{\theta})$ is an estimate of $f(\boldsymbol{\theta})$ obtained by replacing expectation with sample average, and positive definite matrix \mathbf{P} could be a fixed or adaptive preconditioner. By letting $\boldsymbol{\theta}' = \mathbf{P}^{-0.5} \boldsymbol{\theta}$, we can rewrite (2) as

$$\boldsymbol{\theta}' \leftarrow \boldsymbol{\theta}' - \mu \partial \hat{f}(\boldsymbol{\theta}') / \partial \boldsymbol{\theta}' \quad (3)$$

where $\mathbf{P}^{-0.5}$ denotes the principal square root of \mathbf{P} . Hence, (3) suggests that preconditioned SGD is equivalent to SGD in a transformed parameter domain. Within the considered trust region, let us write the stochastic gradient, $\partial \hat{f}(\boldsymbol{\theta}') / \partial \boldsymbol{\theta}'$, explicitly as

$$\partial \hat{f}(\boldsymbol{\theta}') / \partial \boldsymbol{\theta}' = \hat{\mathbf{H}} \boldsymbol{\theta}' + \hat{\mathbf{b}} \quad (4)$$

where $\hat{\mathbf{H}}$ and $\hat{\mathbf{b}}$ are estimates of \mathbf{H} and \mathbf{b} , respectively. Combining (4) and (3) gives the following linear system

$$\boldsymbol{\theta}' \leftarrow (\mathbf{I} - \mu \mathbf{P} \hat{\mathbf{H}}) \boldsymbol{\theta}' - \mu \mathbf{P} \hat{\mathbf{b}} \quad (5)$$

for updating $\boldsymbol{\theta}$ within the assumed trust region, where \mathbf{I} is the identity matrix. A properly determined \mathbf{P} could significantly accelerate convergence of the locally linear system in (5).

We review a few facts shown in (Li, 2018) before introducing our main contributions. Let $\delta \boldsymbol{\theta}$ be a random perturbation of $\boldsymbol{\theta}$, and be small enough such that $\boldsymbol{\theta} + \delta \boldsymbol{\theta}$ still resides in the same trust region. Then, (4) suggests the following resultant perturbation of stochastic gradient,

$$\delta \hat{\mathbf{g}} \stackrel{\text{def}}{=} \partial \hat{f}(\boldsymbol{\theta} + \delta \boldsymbol{\theta}) / \partial \boldsymbol{\theta}' - \partial \hat{f}(\boldsymbol{\theta}') / \partial \boldsymbol{\theta}' = \hat{\mathbf{H}} \delta \boldsymbol{\theta}' = \mathbf{H} \delta \boldsymbol{\theta} + \boldsymbol{\varepsilon} \quad (6)$$

where $\boldsymbol{\varepsilon}$ accounts for the error due to replacing $\hat{\mathbf{H}}$ with \mathbf{H} . Note that by definition, $\delta \hat{\mathbf{g}}$ is a random vector dependent on both \mathbf{z} and $\delta \boldsymbol{\theta}$. The preconditioner in (Li, 2018) is pursued by minimizing criterion

$$c(\mathbf{P}) = E_{z, \delta \boldsymbol{\theta}} [\delta \hat{\mathbf{g}}^T \mathbf{P} \delta \hat{\mathbf{g}} + \delta \boldsymbol{\theta}'^T \mathbf{P}^{-1} \delta \boldsymbol{\theta}'] \quad (7)$$

where subscript $\delta \boldsymbol{\theta}$ in $E_{z, \delta \boldsymbol{\theta}}$ denotes taking expectation over $\delta \boldsymbol{\theta}$. Under mild conditions, criterion (7) determines a unique positive definite \mathbf{P} , which is optimal in the sense that it preconditions the stochastic gradient such that

$$\mathbf{P} E_{z, \delta \boldsymbol{\theta}} [\delta \hat{\mathbf{g}} \delta \hat{\mathbf{g}}^T] \mathbf{P} = E_{\delta \boldsymbol{\theta}} [\delta \boldsymbol{\theta} \delta \boldsymbol{\theta}^T] \quad (8)$$

which is comparable to relationship $\mathbf{H}^{-1}\delta\mathbf{g}\delta\mathbf{g}^T\mathbf{H}^{-1} = \delta\boldsymbol{\theta}\delta\boldsymbol{\theta}^T$, where $\delta\mathbf{g} = \mathbf{H}\delta\boldsymbol{\theta}$ is the perturbation of noiseless gradient, and we assume that \mathbf{H} is invertible, but not necessarily positive definite. Clearly, this preconditioner is comparable to \mathbf{H}^{-1} . It perfectly preconditions the gradient such that the amplitudes of parameter perturbations match that of the associated preconditioned gradient, regardless of the amount of gradient noise. Naturally, preconditioned SGD with this preconditioner inherits the scale-invariance property of the Newton method.

Note that in the presence of gradient noise, the optimal \mathbf{P} and \mathbf{P}^{-1} given by (8) are not unbiased estimates of \mathbf{H}^{-1} and \mathbf{H} , respectively. Actually, even if \mathbf{H} is positive definite and available, \mathbf{H}^{-1} may not always be a good preconditioner when \mathbf{H} is ill-conditioned since it could significantly amplify the gradient noise along the directions of the eigenvectors of \mathbf{H} associated with small eigenvalues, and lead to divergence. More specifically, it is shown in (Li, 2018) that $\mathbf{H}^{-1}E_{z,\delta\theta}[\delta\hat{\mathbf{g}}\delta\hat{\mathbf{g}}^T]\mathbf{H}^{-1} \geq E_{\delta\theta}[\delta\boldsymbol{\theta}\delta\boldsymbol{\theta}^T]$, where $\mathbf{A} \geq \mathbf{B}$ means that $\mathbf{A} - \mathbf{B}$ is nonnegative definite.

3 TWO PRECONDITIONER ESTIMATION CRITERIA

3.1 THE NEWTON TYPE CRITERION

Preconditioner estimation criterion (7) requires $\delta\boldsymbol{\theta}$ to be small enough such that $\boldsymbol{\theta}$ and $\boldsymbol{\theta} + \delta\boldsymbol{\theta}$ reside in the same trust region. In practice, numerical error might be an issue when handling small numbers with floating point arithmetic. This concern becomes more grave with the popularity of single and even half precision math in large scale neural network training. Luckily, (6) relates $\delta\hat{\mathbf{g}}$ to the Hessian-vector product, which can be efficiently evaluated with automatic differentiation software tools. Let \mathbf{v} be a random vector with the same dimension as $\boldsymbol{\theta}$. Then, (4) suggests the following method for Hessian-vector product evaluation,

$$\frac{\partial}{\partial\boldsymbol{\theta}} \left\{ \left[\partial\hat{f}(\boldsymbol{\theta})/\partial\boldsymbol{\theta} \right]^T \mathbf{v} \right\} = \frac{\partial^2\hat{f}(\boldsymbol{\theta})}{\partial\boldsymbol{\theta}\partial\boldsymbol{\theta}^T} \mathbf{v} = \hat{\mathbf{H}}\mathbf{v} \quad (9)$$

Now, replacing $(\delta\boldsymbol{\theta}, \delta\hat{\mathbf{g}})$ in (7) with $(\mathbf{v}, \hat{\mathbf{H}}\mathbf{v})$ leads to our following new preconditioner estimation criterion,

$$c_n(\mathbf{P}) = E_{z,v}[\mathbf{v}^T \hat{\mathbf{H}} \mathbf{P} \hat{\mathbf{H}} \mathbf{v} + \mathbf{v}^T \mathbf{P}^{-1} \mathbf{v}] \quad (10)$$

where the subscript v in $E_{z,v}$ suggests taking expectation over \mathbf{v} . We no longer have the need to assume \mathbf{v} to be an arbitrarily small vector. It is important to note that this criterion only requires the Hessian-vector product. The Hessian itself is not of interest. We call (10) the Newton type preconditioner estimation criterion as the resultant preconditioned SGD method is closely related to the Newton method.

3.2 THE FISHER TYPE CRITERION

We consider the machine learning problems where the *empirical* Fisher information matrix can be well defined by $\mathbf{F} = E_z \left[\frac{\partial\ell(\boldsymbol{\theta}, \mathbf{z})}{\partial\boldsymbol{\theta}} \left(\frac{\partial\ell(\boldsymbol{\theta}, \mathbf{z})}{\partial\boldsymbol{\theta}} \right)^T \right]$. Replacing $(\delta\boldsymbol{\theta}, \delta\hat{\mathbf{g}})$ in (7) with $(\mathbf{v}, \hat{\mathbf{g}} + \lambda\mathbf{v})$ leads to criterion

$$c_f(\mathbf{P}) = E_{z,v}[(\hat{\mathbf{g}} + \lambda\mathbf{v})^T \mathbf{P} (\hat{\mathbf{g}} + \lambda\mathbf{v}) + \mathbf{v}^T \mathbf{P}^{-1} \mathbf{v}] \quad (11)$$

where $\hat{\mathbf{g}} = \partial\hat{f}(\boldsymbol{\theta})/\partial\boldsymbol{\theta}$ is a shorthand for stochastic gradient, and $\lambda \geq 0$ is a damping factor. Clearly, \mathbf{v} is independent of $\hat{\mathbf{g}}$. Let us further assume that \mathbf{v} is drawn from standard multivariate normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$, i.e., $E_v[\mathbf{v}] = \mathbf{0}$ and $E_v[\mathbf{v}\mathbf{v}^T] = \mathbf{I}$. Then, we could simplify $c_f(\mathbf{P})$ as

$$\begin{aligned} c_f(\mathbf{P}) &= \text{tr}\{\mathbf{P}E_{z,v}[(\hat{\mathbf{g}} + \lambda\mathbf{v})(\hat{\mathbf{g}} + \lambda\mathbf{v})^T] + \mathbf{P}^{-1}E_v[\mathbf{v}\mathbf{v}^T]\} \\ &= \text{tr}\{\mathbf{P}[E_z[\hat{\mathbf{g}}\hat{\mathbf{g}}^T] + \lambda^2\mathbf{I}] + \mathbf{P}^{-1}\} \end{aligned} \quad (12)$$

By letting the derivative of $c_f(\mathbf{P})$ with respect to \mathbf{P} be zero, the optimal positive definite solution for $c_f(\mathbf{P})$ is readily shown to be

$$\mathbf{P} = \{E_z[\hat{\mathbf{g}}\hat{\mathbf{g}}^T] + \lambda^2\mathbf{I}\}^{-0.5} \quad (13)$$

When $\hat{\mathbf{g}}$ is a gradient estimation obtained by taking average over B independent samples, $E_z[\hat{\mathbf{g}}\hat{\mathbf{g}}^T]$ is related to the Fisher information matrix by

$$E_z[\hat{\mathbf{g}}\hat{\mathbf{g}}^T] = \mathbf{F}/B + (B-1)\mathbf{g}\mathbf{g}^T/B \quad (14)$$

We call this preconditioner the Fisher type one due to its close relationship to the Fisher information matrix. One can easily modify this preconditioner to obtain an unbiased estimation of \mathbf{F}^{-1} . Let \mathbf{s} be an exponential moving average of $\hat{\mathbf{g}}$. Then, after replacing the $\hat{\mathbf{g}}$ in (13) with $\hat{\mathbf{g}} - \mathbf{s} + \mathbf{s}/\sqrt{B}$ and setting $\lambda = 0$, \mathbf{P}^2/B will be an unbiased estimation of \mathbf{F}^{-1} . Generally, it might be acceptable to keep the bias term, $(B-1)\mathbf{g}\mathbf{g}^T/B$, in (14) for two reasons: it is nonnegative definite and regularizes the inversion in (13); it vanishes when the parameters approach a stationary point. Actually, the Fisher information matrix could be singular for many commonly used models, e.g., finite mixture models, neural networks, hidden Markov models. We might not be able to inverse \mathbf{F} for these singular statistical models without using regularization or damping. A Fisher type preconditioner with $\lambda > 0$ loses the scale-invariance property of a Newton type preconditioner. Both \mathbf{P} and \mathbf{P}^2 can be useful preconditioners when the step size μ and damping factor λ are set properly.

3.3 PROPERTIES OF THE NEWTON TYPE PRECONDITIONER

Following the ideas in (Li, 2018), we can show that (10) determines a unique positive definite preconditioner if and only if $E_v[\mathbf{v}\mathbf{v}^T]$ is positive definite and $\{E_v[\mathbf{v}\mathbf{v}^T]\}^{0.5}E_{z,v}[\hat{\mathbf{H}}\mathbf{v}\mathbf{v}^T\hat{\mathbf{H}}]\{E_v[\mathbf{v}\mathbf{v}^T]\}^{0.5}$ has distinct eigenvalues. Other minimum solutions of criterion (10) are either indefinite or negative definite, and are not interested for our purpose. The proof itself has limited novelty. We omit it here. Instead, let us consider the simplest case, where θ is a scalar parameter, to gain some intuitive understandings of criterion (10). For scalar parameter, it is trivial to show that the optimal solutions minimizing (10) are

$$p = \pm\sqrt{E_v[v^2]/E_{z,v}[\hat{h}^2v^2]} = \pm 1/\sqrt{h^2 + E_z[(h - \hat{h})^2]} \quad (15)$$

where $\hat{\mathbf{H}}, \mathbf{H}, \mathbf{P}$, and \mathbf{v} are replaced with their plain lower case letters, and we have used the fact that $\mathbf{H} - \hat{\mathbf{H}}$ and \mathbf{v} are independent. For gradient descent, we choose the positive solution, although the negative one gives the global minimum of (10). With the positive preconditioner, eigenvalue of the locally linear system in (5) is

$$h/\sqrt{h^2 + E_z[(h - \hat{h})^2]} \quad (16)$$

Now, it is clear that this optimal preconditioner damps the gradient noise when $E_z[(h - \hat{h})^2]$ is large, and preconditions the locally linear system in (5) such that its eigenvalue has unitary amplitude when the gradient noise vanishes. Convergence is ensured when a normalized step size, i.e., $0 < \mu < 1$, is used. For θ with higher dimensions, eigenvalues of the locally linear system in (5) is normalized into range $[-1, 1]$ as well, in a way similar to (16).

4 ESTIMATING PRECONDITIONERS ON MATRIX LIE GROUPS

4.1 UPDATING RULE FOR THE NEWTON TYPE PRECONDITIONER

Let us take the Newton type preconditioner as an example to derive its updating rule. Updating rule for the Fisher type preconditioner is the same except for replacing the Hessian-vector product with stochastic gradient. Here, Lie group always refers to the matrix Lie group.

It is inconvenient to optimize \mathbf{P} directly as it must be a symmetric and positive definite matrix. Instead, we represent the preconditioner as $\mathbf{P} = \mathbf{Q}^T\mathbf{Q}$, and estimate \mathbf{Q} . Now, \mathbf{Q} must be a nonsingular matrix as both $c_n(\mathbf{P})$ and $c_f(\mathbf{P})$ diverge when \mathbf{P} is singular. Invertible matrices with the same dimension form a Lie group. In practice, we are more interested in Lie groups with sparse representations. Examples of such groups are given in the next section. Let us consider a proper small perturbation of \mathbf{Q} , $\delta\mathbf{Q}$, such that $\mathbf{Q} + \delta\mathbf{Q}$ still lives on the same Lie group. The distance between \mathbf{Q} and $\mathbf{Q} + \delta\mathbf{Q}$ can be naturally defined as $\text{dist}(\mathbf{Q}, \mathbf{Q} + \delta\mathbf{Q}) = \sqrt{\text{tr}(\delta\mathbf{Q}\mathbf{Q}^{-1}\mathbf{Q}^{-T}\delta\mathbf{Q}^T)}$ (Amari, 1998). Intuitively, this distance is larger for the same amount of perturbation when \mathbf{Q} is closer to a singular matrix. With the above tensor metric definition, natural gradient for optimizing \mathbf{Q} has form

$$\nabla\mathbf{Q} = \mathbf{R}\mathbf{Q} \quad (17)$$

For example, when \mathbf{Q} lives on the group of invertible upper triangular matrices, \mathbf{R} is given by

$$\mathbf{R} = 2\text{triu}\{E_{z,v}[Q\hat{\mathbf{H}}\mathbf{v}\mathbf{v}^T\hat{\mathbf{H}}^T\mathbf{Q}^T - \mathbf{Q}^{-T}\mathbf{v}\mathbf{v}^T\mathbf{Q}^{-1}]\} \quad (18)$$

where triu takes the upper triangular part of a matrix. Another way to derive (17) is to let $\delta\mathbf{Q} = \mathcal{E}\mathbf{Q}$, and consider the derivative with respect to \mathcal{E} , where \mathcal{E} is a proper small matrix such that $\mathbf{Q} + \mathcal{E}\mathbf{Q}$ still lives on the same Lie group. Gradient derived in this way is known as relative gradient (Cardoso & Laheld, 1996). For our preconditioner learning problem, relative gradient and natural gradient have the same form. Now, \mathbf{Q} can be updated using natural or relative gradient descent as

$$\mathbf{Q} \leftarrow \mathbf{Q} - \mu_q \mathbf{R}\mathbf{Q} \quad (19)$$

In practice, it is convenient to use the following updating rule with normalized step size,

$$\mathbf{Q} \leftarrow (\mathbf{I} - \mu_0 \mathbf{R}/\|\mathbf{R}\|)\mathbf{Q} \quad (20)$$

where $0 < \mu_0 < 1$, and $\|\cdot\|$ takes the norm of a matrix. One simple choice for matrix norm is the maximum absolute value of a matrix.

Note that natural gradient can take different forms. One should not confuse the natural gradient on the Lie group derived from a tensor metric with the natural gradient for model parameter learning derived from a Fisher information metric.

4.2 SUMMARY OF THE NEWTON TYPE PRECONDITIONED SGD

One iteration of the Newton type preconditioned SGD consists of the following steps.

1. Evaluate stochastic gradient $\hat{\mathbf{g}}$.
2. Draw \mathbf{v} from $\mathcal{N}(\mathbf{0}, \mathbf{I})$, and evaluate Hessian-vector product $\hat{\mathbf{H}}\mathbf{v}$.
3. Update preconditioners with $\mathbf{Q} \leftarrow (\mathbf{I} - \mu_0 \hat{\mathbf{R}}/\|\hat{\mathbf{R}}\|)\mathbf{Q}$.
4. Update parameters with $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \mu\mathbf{Q}^T\mathbf{Q}\hat{\mathbf{g}}$.

The two step sizes, μ and μ_0 , are normalized. They should take values in range $[0, 1]$ with typical value 0.01. We usually initialize \mathbf{Q} to a scaled identity matrix with proper dimension. The specific form of $\hat{\mathbf{R}}$ depends on the Lie group to be considered. For example, for upper triangular \mathbf{Q} , we have $\hat{\mathbf{R}} = 2\text{triu}[Q\hat{\mathbf{H}}\mathbf{v}\mathbf{v}^T\hat{\mathbf{H}}^T\mathbf{Q}^T - \mathbf{Q}^{-T}\mathbf{v}\mathbf{v}^T\mathbf{Q}^{-1}]$, where $\mathbf{Q}^{-T}\mathbf{v}$ can be efficiently calculated with back substitution.

4.3 SUMMARY OF THE FISHER TYPE PRECONDITIONED SGD

We only need to replace $\hat{\mathbf{H}}\mathbf{v}$ in the Newton type preconditioned SGD with $\hat{\mathbf{g}} + \lambda\mathbf{v}$ to obtain the Fisher type one. Thus, we do not list its main steps here. Note that only its step size for the preconditioner updating is normalized. There is no simple way to jointly determine the proper ranges for step size μ and damping factor λ . Again, $\hat{\mathbf{R}}$ may take different forms on different Lie groups. For upper triangular \mathbf{Q} , we have $\hat{\mathbf{R}} = 2\text{triu}[Q(\hat{\mathbf{g}} + \lambda\mathbf{v})(\hat{\mathbf{g}} + \lambda\mathbf{v})^T\mathbf{Q}^T - \mathbf{Q}^{-T}\mathbf{v}\mathbf{v}^T\mathbf{Q}^{-1}]$, where $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Here, it is important to note that the natural or relative gradient for $c_f(\mathbf{P})$ with the form given in (12) involves explicit matrix inversion. However, matrix inversion can be avoided by using the $c_f(\mathbf{P})$ in (11), which includes \mathbf{v} as an auxiliary variable. It is highly recommended to avoid explicit matrix inversion for large \mathbf{Q} .

4.4 VARIATIONS OF PRECONDITIONED SGD

There are many ways to modify the above preconditioned SGD methods. Since curvatures typically evolves slower than gradients, one can update the preconditioner less frequently to save computations per iteration. With parallel computing available, one might update the preconditioner and model parameters simultaneously and asynchronously to save wall time per iteration. Combining preconditioner and momentum may further accelerate convergence. For recurrent neural network learning, we may need to clip the norm of preconditioned gradients to avoid excessively large parameter updates. In general, preconditioned gradient clipping relates to the trust region method by

$$\|\boldsymbol{\theta}^{\text{[new]}} - \boldsymbol{\theta}\| = \|\mu\mathbf{P}\hat{\mathbf{g}}/\max(1, \|\mathbf{P}\hat{\mathbf{g}}\|/\Omega)\| \leq \mu\Omega$$

where $\Omega > 0$ is a clipping threshold, comparable to the size of trust region. One may set Ω to a positive number proportional to the square root of the number of model parameters. Most importantly, we can choose different Lie groups for estimating our preconditioners to achieve a good trade off between performance and complexity.

5 USEFUL LIE GROUPS WITH SPARSE REPRESENTATIONS

In practice, we seldom consider the Lie group consisting of dense invertible matrices for preconditioner estimation when the problem size is large. Lie groups with sparse structures are of more interests. To begin with, let us recall a few facts about Lie group. If \mathbf{A} and \mathbf{B} are two Lie groups, then \mathbf{A}^T , $\mathbf{A} \otimes \mathbf{B}$, and $\mathbf{A} \oplus \mathbf{B}$ all are Lie groups, where \otimes and \oplus denote Kronecker product and direct sum, respectively. Furthermore, for any matrix \mathbf{C} with compatible dimensions, block matrix

$$\mathbf{Q} = \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{0} & \mathbf{B} \end{bmatrix} \quad (21)$$

still forms a Lie group. We do not show proofs of the above statements here as they are no more than a few lines of algebraic operations. These simple rules can be used to design many useful Lie groups for constructing sparse preconditioners. We already know that invertible upper triangular matrices form a Lie group. Here, we list a few useful ones with sparse representations.

5.1 DIAGONAL PRECONDITIONER

Diagonal matrices with the same dimension and positive diagonal entries form a Lie group with reducible representation. Preconditioners learned on this group are called diagonal preconditioners.

5.2 KRONECKER PRODUCT PRECONDITIONER

For matrix parameter Θ , we can flatten Θ into a vector, and precondition its gradient using a Kronecker product preconditioner with \mathbf{Q} having form $\mathbf{Q} = \mathbf{Q}_2 \otimes \mathbf{Q}_1$. Clearly, \mathbf{Q} is a Lie group as long as \mathbf{Q}_1 and \mathbf{Q}_2 are two Lie groups. Let us check its role in learning the following affine transformation

$$\mathbf{y} = [\Theta_{\text{weight}}, \Theta_{\text{bias}}] \mathbf{x} = \Theta \mathbf{x} \quad (22)$$

where \mathbf{x} is the input feature vector augmented with 1, and \mathbf{y} is the output feature vector. After reverting the flattened Θ back to its matrix form, the preconditioned SGD learning rule for Θ is

$$\Theta \leftarrow \Theta - \mu \mathbf{Q}_1^T \mathbf{Q}_1 \frac{\partial \hat{f}(\Theta)}{\partial \Theta} \mathbf{Q}_2^T \mathbf{Q}_2 \quad (23)$$

Similar to (3), we introduce coordinate transformation $\Theta' = \mathbf{Q}_1^{-T} \Theta \mathbf{Q}_2^{-1}$, and rewrite (23) as

$$\Theta' \leftarrow \Theta' - \mu \partial \hat{f}(\Theta) / \partial \Theta' \quad (24)$$

Correspondingly, the affine transformation in (22) is rewritten as $\mathbf{y}' = \Theta' \mathbf{x}'$, where $\mathbf{y}' = \mathbf{Q}_1^{-T} \mathbf{y}$ and $\mathbf{x}' = \mathbf{Q}_2 \mathbf{x}$ are the transformed input and output feature vectors, respectively. Hence, the preconditioned SGD in (23) is equivalent to the SGD in (24) with transformed feature vectors \mathbf{x}' and \mathbf{y}' . We know that feature whitening and normalization could significantly accelerate convergence. A Kronecker product preconditioner plays a similar role in learning the affine transformation in (22).

5.3 SCALING AND NORMALIZATION PRECONDITIONER

This is a special Kronecker product preconditioner by constraining \mathbf{Q}_1 to be a diagonal matrix, and \mathbf{Q}_2 to be a sparse matrix where only its diagonal and last column can have nonzero values. Note that \mathbf{Q}_2 with nonzero diagonal entries forms a Lie group. Hence, $\mathbf{Q} = \mathbf{Q}_2 \otimes \mathbf{Q}_1$ is a Lie group as well. We call it a scaling and normalization preconditioner as it resembles a preconditioner that scales the output features and normalizes the input features. Let us check the transformed features $\mathbf{y}' = \mathbf{Q}_1^{-T} \mathbf{y}$ and $\mathbf{x}' = \mathbf{Q}_2 \mathbf{x}$. It is clear that \mathbf{y}' is an element-wisely scaled version of \mathbf{y} as \mathbf{Q}_1 is a diagonal matrix. To make \mathbf{x}' a ‘‘normalized’’ feature vector, \mathbf{x} needs to be an input feature vector augmented with 1.

Let us check a simple example to verify this point. We consider an input vector with two features, and write down its normalized features explicitly as below,

$$\begin{bmatrix} x'_1 \\ x'_2 \\ 1 \end{bmatrix} = \begin{bmatrix} 1/\sigma_1 & 0 & -m_1/\sigma_1 \\ 0 & 1/\sigma_2 & -m_2/\sigma_2 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix} \quad (25)$$

where m_i and σ_i are the mean and standard deviation of x_i , respectively. It is straightforward to show that the feature normalization operation in (25) forms a Lie group with four freedoms. For the scaling-and-normalization preconditioner, we have no need to force the last diagonal entry of \mathbf{Q}_2 to be 1. Hence, the group of feature normalization operation is a subgroup of \mathbf{Q}_2 .

5.4 SCALING AND WHITENING PRECONDITIONER

This is another special Kronecker product preconditioner by constraining \mathbf{Q}_1 to be a diagonal matrix, and \mathbf{Q}_2 to be an upper triangular matrix with positive diagonal entries. We call it a scaling-and-whitening preconditioner since it resembles a preconditioner that scales the output features and whitens the input features. Again, the input feature vector \mathbf{x} must be augmented with 1 such that the whitening operation forms a Lie group represented by upper triangular matrices with 1 being its last diagonal entry. This is a subgroup of \mathbf{Q}_2 as we have no need to fix \mathbf{Q}_2 's last diagonal entry to 1.

It is not possible to enumerate all kinds of Lie groups suitable for constructing preconditioners. For example, Kronecker product preconditioner with form $\mathbf{Q} = \mathbf{Q}_3 \otimes \mathbf{Q}_2 \otimes \mathbf{Q}_1$ could be used for preconditioning gradients of a third order tensor. The normalization and whitening groups are just two special cases of the groups with the form shown in (21), and there are numerous more choices having sparsities between that of these two. Regardless of the detailed form of \mathbf{Q} , all such preconditioners share the same form of learning rule shown in (20), and they all can be efficiently learned using natural or relative gradient descent without much tuning effort.

6 RELATIONSHIP TO EXISTING METHODS

6.1 RELATIONSHIP TO ADAGRAD, RMSPROP AND ADAM

Adagrad, RMSProp and Adam all use the Fisher type preconditioner living on the group of diagonal matrices with positive diagonal entries. This is a simple group. Optimal solution for $c_f(\mathbf{P})$ has closed-form solution $\mathbf{P} = \text{diag}(1 \oslash \sqrt{E_z[\hat{\mathbf{g}} \odot \hat{\mathbf{g}}] + \lambda^2})$, where \odot and \oslash denote element wise multiplication and division, respectively. In practice, simple exponential moving average is used to replace the expectation when using this preconditioner.

6.2 RELATIONSHIP TO EQUILIBRATED SGD

For diagonal preconditioner, the optimal solution minimizing $c_n(\mathbf{P})$ has closed-form solution $\mathbf{P} = \text{diag}\left(\sqrt{E_v[\mathbf{v} \odot \mathbf{v}] \oslash E_{z,v}[\hat{\mathbf{H}}\mathbf{v} \odot \hat{\mathbf{H}}\mathbf{v}]}\right)$. For $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $E_v[\mathbf{v} \odot \mathbf{v}]$ reduces to a vector with unit entries. Then, this optimal solution gives the equilibration preconditioner in (Dauphin et al., 2015).

6.3 RELATIONSHIP TO KFAC AND SIMILAR METHODS

The preconditioners considered in (Povey et al., 2015) and (Martens & Grosse, 2015) are closely related to the Fisher type Kronecker product preconditioners. While KFAC approximates the Fisher metric of a matrix parameter as a Kronecker product to obtain its approximated inverse in closed-form solution, our method turns to an iterative solution to approximate this same inverse. Theoretically, our method's accuracy is only limited by the expressive power of the Lie group since no intermediate approximation is made. In practice, one distinct advantage of our method over KFAC is that explicit matrix inversion is avoided by introducing auxiliary vector \mathbf{v} and using back substitution, while KFAC typically requires inversion of symmetric matrices. On graphics processing units (GPU) and with large matrices, parallel back substitution is as computationally cheap as matrix multiplication, and could be several orders of magnitude faster than inversion of symmetric matrix. Another advantage is that our method is derived from a unified framework. There is no need to invent different preconditioner learning rules when we switch the Lie group representations.

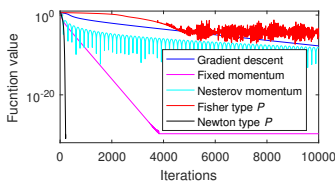


Figure 1: Convergence curves of compared methods on minimizing the Rosenbrock function.

6.4 RELATIONSHIP TO BATCH NORMALIZATION

Batch normalization can be viewed as preconditioned SGD using a specific scaling-and-normalization preconditioner with constraint $\mathbf{Q}_1 = \mathbf{I}$ and \mathbf{Q}_2 from the feature normalization Lie group. However, we should be aware that explicit input feature normalization is only empirically shown to accelerate convergence, and has little meaning in certain scenarios, e.g., recurrent neural network learning where features may not have any stationary first or second order statistic. Both the Newton and Fisher type preconditioned SGD methods provide a more general and principled approach to find the optimal preconditioner, and apply to a broader range of applications. Generally, a scaling-and-normalization preconditioner does not necessarily “normalize” the input features in the sense of mean removal and variance normalization.

7 EXPERIMENTAL RESULTS

We use the square root Fisher type preconditioners in the following experiments since they are less picky on the damping factor, and seem to be more numerically robust on large scale problems. Still, as shown in our Pytorch implementation package, the original Fisher type preconditioners could perform better on small scale problems like the MNIST handwritten digit recognition task.

7.1 APPLICATION TO MATHEMATICAL OPTIMIZATION

Let us consider the minimization of Rosenbrock function, $f(\boldsymbol{\theta}) = 100(\theta_2 - \theta_1^2)^2 + (1 - \theta_1)^2$, starting from initial guess $\boldsymbol{\theta} = [-1, 1]$. This is a well known benchmark problem for mathematical optimization. The compared methods use fixed step size. For each method, the best step size is selected from sequence $\{\dots, 1, 0.5, 0.2, 0.1, 0.05, 0.02, 0.01, \dots\}$. For gradient descent, the best step size is 0.002. For momentum method, the moving average factor is 0.9, and the best step size is 0.002. For Nesterov momentum, the best step size is 0.001. For preconditioned SGD, \mathbf{Q} is initialized to $0.1\mathbf{I}$ and lives on the group of triangular matrices. For the Fisher type method, we set $\lambda = 0.1$, and step sizes 0.01 and 0.001 for preconditioner and parameter updates, respectively. For the Newton type method, we set step sizes 0.2 and 0.5 for preconditioner and parameter updates, respectively. Figure 1 summarizes the results. The Newton type method performs the best, converging to the optimal solution using about 200 iterations. The Fisher type method does not fit into this problem, and performs poorly as expected. Mathematical optimization is not our focus. Still, this example shows that the Newton type preconditioned SGD works well for mathematical optimization.

7.2 IMAGENET EXPERIMENT

We consider the ImageNet ILSVRC2012 database for the image classification task. The well known AlexNet is considered. We follow the descriptions in (Alex et al., 2012) as closely as possible to set up our experiment. One main difference is that we do not augment the training data. Another big difference is that we use a modified local response normalization (LRN). The LRN function from TensorFlow implementation is not second order differentiable. We have to approximate the local energy used for LRN with a properly scaled global energy to facilitate Hessian-vector product evaluation. Note that convolution can be rewritten as correlation between the flattened input image patches and filter coefficients. In this way, we find that there are eight matrices to be optimized in the AlexNet, and their shapes are: $[96, 11 \times 11 \times 3 + 1]$, $[256, 5 \times 5 \times 96 + 1]$, $[384, 3 \times 3 \times 256 + 1]$, $[384, 3 \times 3 \times 384 + 1]$, $[256, 3 \times 3 \times 384]$, $[4096, 6 \times 6 \times 256 + 1]$, $[4096, 4096 + 1]$, and $[1000, 4096 + 1]$. We have tried diagonal and scaling-and-normalization preconditioners for

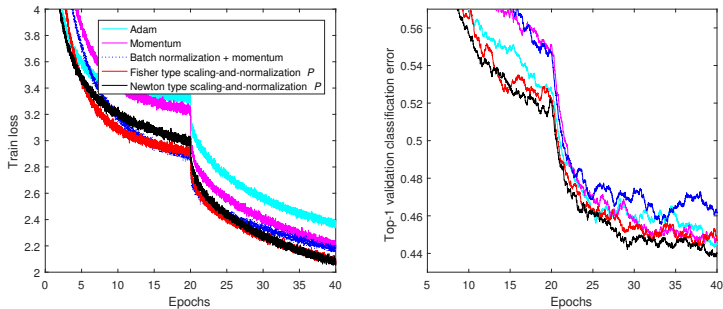


Figure 2: Typical smoothed learning curves from compared methods for the ImageNet ILSVRC2012 image classification task with AlexNet. Note that batch normalization alters the L2-regularization. Its training loss is not directly comparable with others.

each matrix. Denser preconditioners, e.g., the Kronecker product one, require hundreds of millions parameters for representations, and are expensive to run on our platform. Each compared method is trained with 40 epochs, mini-batch size 128, step size μ for the first 20 epochs, and 0.1μ for the last 20 epochs. We have compared several methods with multiple settings, and only report the ones with reasonably good results here. For Adam, the initial step size is set to 0.00005. For batch normalization, initial step size is 0.002, and its moving average factors for momentum and statistics used for feature normalization are 0.9 and 0.99, respectively. The momentum method uses initial step size 0.002, and moving average factor 0.9 for momentum. Preconditioned SGD performs better with the scaling-and-normalization preconditioner. Its \mathbf{Q} is initialized to $0.1\mathbf{I}$, and updated with normalized step size 0.01. For the Fisher type preconditioner, we set $\lambda = 0.001$ and initial step size 0.00005. For the Newton type preconditioner, its initial step size is 0.01. Figure 2 summarizes the results. Training loss for batch normalization is only for reference purpose as normalization alters the L2-regularization term. Batch normalization does not perform well under this setup, maybe due to its conflict with certain settings like the LRN and L2-regularization. We see that the scaling-and-normalization preconditioner does accelerate convergence, although it is super sparse. The Newton type preconditioned SGD performs the best, and achieves top-1 validation accuracy about 56% when using only one crop for testing, while the momentum method may require 90 epochs to achieve similar performance.

7.3 WORD LEVEL LANGUAGE MODELING EXPERIMENT

We consider the world level language modeling problem with reference implementation available from <https://github.com/pytorch/examples>. The Wikitext-2 database with 33278 tokens is considered. The task is to predict the next token from history observations. Our tested network consists of six layers, i.e., encoding layer, LSTM layer, dropout layer, LSTM layer, dropout layer, and decoding layer. For each LSTM layer, we put all its coefficients into a single matrix Θ by defining output and augmented input feature vectors as in $[\hat{\mathbf{i}}_t; \mathbf{f}_t; \mathbf{g}_t; \mathbf{o}_t] = \Theta [\mathbf{x}_t; \mathbf{h}_{t-1}; 1]$, $\mathbf{c}_t = \mathbf{f}_t \mathbf{c}_{t-1} + \hat{\mathbf{i}}_t \mathbf{g}_t$, $\mathbf{h}_t = \mathbf{o}_t \tanh(\mathbf{c}_t)$, where t is a discrete time index, \mathbf{x} is the input, \mathbf{h} is the hidden state, and \mathbf{c} is the cell state. The encoding layer’s weight matrix is the transpose of that of the decoding layer. Thus, we totally get three matrices to be optimized. With hidden layer size 200, shapes of these three matrices are $[4 \times 200, 2 \times 200 + 1]$, $[4 \times 200, 2 \times 200 + 1]$, and $[33278, 200 + 1]$, respectively. For all methods, the step size is reduced to one fourth of the current value whenever the current perplexity on validation set is larger than the best one ever found. For SGD, the initial step size is 20, and the gradient is clipped with threshold 0.25. The momentum method diverges easily without clipping. We set momentum 0.9, initial step size 1, and clipping threshold 0.25. We set initial step size 0.005 and damping factor $\lambda^2 = 10^{-12}$ for Adam and sparse Adam. Sparse Adam updates its moments and model parameters only when the corresponding stochastic gradients are not zeros. We have tried diagonal, scaling-and-normalization and scaling-and-whitening preconditioners for each matrix. The encoding (decoding) matrix is too large to consider KFAC like preconditioner. The diagonal preconditioner performs the worst, and the other two have comparable performance. For both types of preconditioned SGD, the clipping threshold for preconditioned gradient is 100, the initial step size is 0.1, and \mathbf{Q} is initialized to \mathbf{I} . We set $\lambda = 0$ for the Fisher

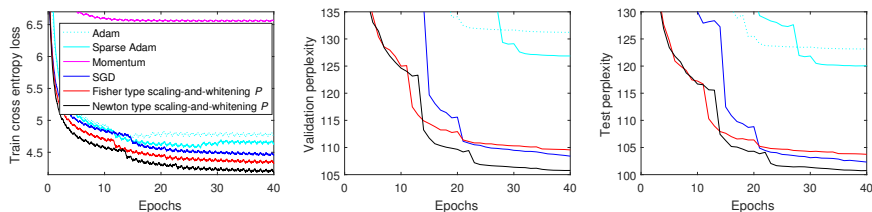


Figure 3: Typical learning curves from compared methods for the Wikitext-2 database word level language modeling task with a LSTM neural network.

type preconditioned SGD. With dropout rate 0.2, the best three test perplexities are 100.7 by SGD, 103.9 by Newton type preconditioned SGD, and 105.8 by Fisher type preconditioned SGD. With dropout rate 0.35, the best three test perplexities are 100.7 by Newton type preconditioned SGD, 102.3 by SGD, and 103.7 by Fisher type preconditioned SGD. Although SGD performs well on the test set, both types of preconditioned SGD have significantly lower training losses than SGD with either dropout rate. Figure 3 summarizes the results when the dropout rate is 0.35. Methods involving momentum, including Adam and sparse Adam, perform poorly. Note that our preconditioners preserve the sparsity property of gradients from the encoding and decoding layers (Appendix A). This saves considerable computations by avoiding update parameters with zero gradients. Again, both preconditioners accelerate convergence significantly despite their high sparsity.

7.4 COMPUTATIONAL COMPLEXITY AND IMPLEMENTATION

Compared with SGD, the Fisher type preconditioned SGD adds limited computational complexity when sparse preconditioners are adopted. The Newton type preconditioned SGD requires Hessian-vector product, which typically has complexity comparable to that of gradient evaluation. Thus, using SGD as the base line, the Newton type preconditioned SGD approximately doubles the computational complexity per iteration, while the Fisher type SGD has similar complexity. Wall time per iteration of preconditioned SGD highly depends on the implementations. Ideally, the preconditioners and parameters could be updated in a parallel and asynchronous way such that SGD and preconditioned SGD have comparable wall time per iteration.

We have put our TensorFlow and Pytorch implementations on <https://github.com/lixilinx>. More experimental results comparing different preconditioners and optimization methods on diverse benchmark problems can be found there. For the ImageNet experiment, all compared methods are implemented in Tensorflow, and require two days and a few hours to finish 40 epochs on a GeForce GTX 1080 Ti GPU. The word level language modeling experiment is implemented in Pytorch. We have rewritten the word embedding function to enable second order derivative. For this task, SGD and the Fisher type preconditioned SGD have similar wall time per iteration, while the Newton type method requires about 80% more wall time per iteration than SGD when running on the same GPU.

8 CONCLUSIONS

Two types of preconditioners and preconditioned SGD methods are studied. The one requiring Hessian-vector product for preconditioner estimation is suitable for general purpose optimization. We call it the Newton type preconditioned SGD due to its close relationship to the Newton method. The other one only requires gradient for preconditioner estimation. We call it the Fisher type preconditioned SGD as its preconditioner is closely related to the inverse of Fisher information matrix. Both preconditioners can be efficiently learned using natural or relative gradient descent on any matrix Lie groups designated by the user. The Fisher type preconditioned SGD has lower computational complexity per iteration, but may require more tuning efforts on selecting its step size and damping factor. The Newton type preconditioned SGD has higher computational complexity per iteration, but is more user friendly due to its use of normalized step size and built-in gradient noise damping ability. Both preconditioners, even with very sparse representations, are shown to considerably accelerate convergence on relatively large scale problems.

ACKNOWLEDGMENTS

The author thanks the reviewers for their comments and Yaroslav Bulatov for his discussions to improve this paper.

REFERENCES

- K. Alex, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *NIPS*, pp. 1097–1105, 2012.
- S. Amari. Natural gradient works efficiently in learning. *Neural Computation*, 10(2):251–276, 1998.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- J. F. Cardoso and B. H. Laheld. Equivariant adaptive source separation. *IEEE Trans. Signal Process.*, 44(12):3017–3030, 1996.
- Y. N. Dauphin, H. Vries, and Y. Bengio. Equilibrated adaptive learning rates for non-convex optimization. In *NIPS*, pp. 1504–1512. MIT Press, 2015.
- D. John, H. Elad, and S. Yoram. Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159, 2011.
- D. P. Kingma and J. L. Ba. Adam: a method for stochastic optimization. In *ICLR*. Ithaca, NY: arXiv.org, 2015.
- X. L. Li. Preconditioned stochastic gradient descent. *IEEE Trans. Neural Networks and Learning Systems*, 29(5):1454–1466, 2018.
- J. Martens and R. B. Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In *ICML*, pp. 2408–2417, 2015.
- Y. Nesterov. A method of solving a convex programming problem with convergence rate $o(1/\sqrt{k})$. *Soviet Mathematics Doklady*, 27:372–376, 1983.
- D. Povey, X. Zhang, and S. Khudanpur. Parallel training of DNNs with natural gradient and parameter averaging. In *ICLR*. Ithaca, NY: arXiv.org, 2015.
- D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors. *Nature*, 323:533–536, 1986.

APPENDIX A: WHITENING-AND-SCALING PRECONDITIONER PRESERVES SPARSITY OF GRADIENT OF ENCODING MATRIX

We are to show that the scaling-and-whitening and whitening-and-scaling preconditioners preserve the sparsity property of gradients from the decoding and encoding layers in the word level language model considered in Experiment 3, respectively. Clearly, we only need to show this for the encoding part. Let

$$\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_I]$$

be a $D \times I$ word embedding matrix, where D is the embedding dimension, I is the number of tokens, and \mathbf{w}_i is the vector representation for the i th token. Typically, only a whitening-and-scaling and sparser preconditioners are affordable since $I \gg D > 1$. Notably, for sufficiently large I , the whitening-and-scaling preconditioner could be sparser than the diagonal one. Let us consider preconditioner

$$\mathbf{P} = \mathbf{Q}^T \mathbf{Q}, \quad \mathbf{Q} = \mathbf{Q}_2 \otimes \mathbf{Q}_1, \quad \mathbf{Q}_2 = \text{diag}(q_{2,1}, q_{2,2}, \dots, q_{2,I})$$

By (23), the preconditioned gradient is

$$[q_{2,1}^2 \mathbf{Q}_1^T \mathbf{Q}_1 \hat{\mathbf{g}}_1, q_{2,2}^2 \mathbf{Q}_1^T \mathbf{Q}_1 \hat{\mathbf{g}}_2, \dots, q_{2,I}^2 \mathbf{Q}_1^T \mathbf{Q}_1 \hat{\mathbf{g}}_I]$$

where $\hat{\mathbf{g}}_i$ is the stochastic gradient for the i th word embedding vector. Since $I \gg B$, most of these $\hat{\mathbf{g}}_i$'s are zeros, where B is the batch size. This preconditioner does not mix up gradients of different word embedding vectors. Hence, the sparsity property is preserved.