

Appendix

We analyze the geometric structure of VPNG to show its insight. Then we provide more details for the experiments.

A. Analysis on the geometric structure of VPNG

As discussed in Hoffman et al. (2013), the traditional natural gradient points to the steepest ascent direction of the ELBO in the symmetric KL divergence space of the variational distribution q . Mathematically, for the ELBO function as in Equation 1, the traditional natural gradient points to the direction of the solution to the following optimization problem, as $\epsilon \rightarrow 0$:

$$\begin{aligned} & \arg \max_{\Delta \lambda} \mathcal{L}(\lambda + \Delta \lambda, \theta) \\ & \text{s.t. } \text{KL}_{\text{sym}}(q(\lambda) \| q(\lambda + \Delta \lambda)) \leq \epsilon. \end{aligned}$$

In fact, denote $\eta = \begin{pmatrix} \lambda \\ \theta \end{pmatrix}$, the reparameterization for $\mathbf{z} = g(\mathbf{x}, \varepsilon; \lambda)$ and $p_{\mathbf{x}'}(\eta) = p(\mathbf{x}' | \mathbf{z} = g(\mathbf{x}, \varepsilon; \lambda); \theta)$ to be the reparameterized predictive distribution, our VPNG (as defined in Equation 14) shares similar geometric structures and points to the direction of the solution to the following optimization problem, as $\epsilon \rightarrow 0$:

$$\begin{aligned} & \arg \max_{\Delta \eta} \mathcal{L}(\eta + \Delta \eta) \\ & \text{s.t. } \mathbb{E}_{\varepsilon} [\text{KL}_{\text{sym}}(p_{\mathbf{x}'}(\eta) \| p_{\mathbf{x}'}(\eta + \Delta \eta))] \leq \epsilon. \end{aligned} \quad (18)$$

Here the expectation on ε takes with respect to the parameter-free distribution $s(\varepsilon)$ in the reparameterization.

Proof. The proof for the above fact is similar with the proof for the traditional natural gradient as in Hoffman et al. (2013). Ideally, we want to find a (possibly approximate) Riemannian metric $G(\eta)$ to capture the geometric structure of the expected symmetric KL divergence $\mathbb{E}_{\varepsilon} [\text{KL}_{\text{sym}}(p_{\mathbf{x}'}(\eta) \| p_{\mathbf{x}'}(\eta + \Delta \eta))]$:

$$\begin{aligned} \mathbb{E}_{\varepsilon} [\text{KL}_{\text{sym}}(p_{\mathbf{x}'}(\eta) \| p_{\mathbf{x}'}(\eta + \Delta \eta))] & \approx \Delta \eta^{\top} G(\eta) \Delta \eta \\ & + o(\|\Delta \eta\|^2). \end{aligned}$$

By making first-order Taylor approximation on $p_{\mathbf{x}'}(\eta + \Delta \eta)$

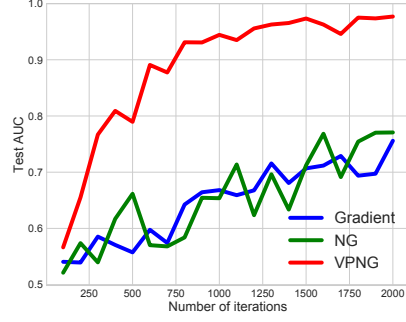


Figure 4. Bayesian Logistic regression test AUC-iteration learning curve.

and $\log p_{\mathbf{x}'}(\eta + \Delta \eta)$, we get

$$\begin{aligned} & \mathbb{E}_{\varepsilon} [\text{KL}_{\text{sym}}(p_{\mathbf{x}'}(\eta) \| p_{\mathbf{x}'}(\eta + \Delta \eta))] \\ & = \mathbb{E}_{\varepsilon} \left[\int (p_{\mathbf{x}'}(\eta + \Delta \eta) - p_{\mathbf{x}'}(\eta)) \right. \\ & \quad \left. \cdot (\log p_{\mathbf{x}'}(\eta + \Delta \eta) - \log p_{\mathbf{x}'}(\eta)) dx' \right] \\ & = \mathbb{E}_{\varepsilon} \left[\int (\nabla_{\eta} p_{\mathbf{x}'}(\eta)^{\top} \Delta \eta) \cdot (\nabla_{\eta} \log p_{\mathbf{x}'}(\eta)^{\top} \Delta \eta) dx' \right. \\ & \quad \left. + O(\|\Delta \eta\|^3) \right] \\ & = \mathbb{E}_{\varepsilon} \left[\int p_{\mathbf{x}'}(\eta) \cdot (\nabla_{\eta} \log p_{\mathbf{x}'}(\eta)^{\top} \Delta \eta) \right. \\ & \quad \left. \cdot (\nabla_{\eta} \log p_{\mathbf{x}'}(\eta)^{\top} \Delta \eta) dx' + O(\|\Delta \eta\|^3) \right] \\ & = \Delta \eta^{\top} F_r \Delta \eta + O(\|\Delta \eta\|^3). \end{aligned}$$

The term $O(\|\Delta \eta\|^3)$ is negligible compared to the first term when $\epsilon \rightarrow 0$. Hence, we could take $G(\eta)$ to be just F_r , the variational predictive Fisher information as defined in Section 3.3. By Amari (1998)'s analysis on natural gradients, we know that the solution to Equation 18 points to the direction of $G(\eta)^{-1} \cdot \nabla_{\lambda, \theta} \mathcal{L} = \nabla_{\lambda, \theta}^{\text{VPNG}} \mathcal{L}$, when $\epsilon \rightarrow 0$. \square

B. More details for the experiments

For the Bayesian Logistic regression experiment, we show the test AUC-iteration curve as in Figure 4. It can be seen that the VPNG behaves more stable compared to the baseline methods.

For the VAE and the VMF experiments, we chose hyperparameters for all methods based on the training ELBO at the end of the time budget.

For the traditional natural gradient and VPNG, we applied the damping factor μ . More precisely, we take VPNG updates as $\hat{\nabla}_{\theta, \lambda}^{\text{VPNG}} \mathcal{L} = (\hat{F}_r + \mu I)^{-1} \cdot \nabla_{\theta, \lambda} \mathcal{L}$ and traditional natural gradient updates as $\hat{\nabla}_{\lambda}^{\text{NG}} \mathcal{L} = (\hat{F}_q + \mu I)^{-1} \cdot \nabla_{\lambda} \mathcal{L}$. This is also applied in the Bayesian Logistic regression experiment in Section 5.1.

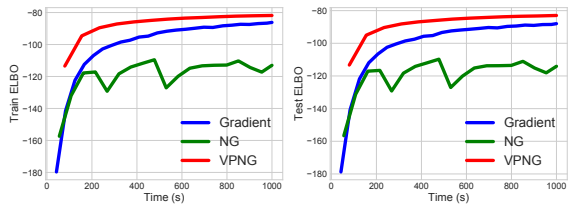


Figure 5. VAE learning curves on binarized MNIST, without exponential moving averages

For the VAE and the VMF experiments, we applied the K-FAC approximation (Martens & Grosse, 2015) to efficiently approximate the Fisher information matrices, in the NG and VPNG computations. For the VPNGs, we view the VAE model as a 6-layer neural network and the VMF model as a 4-layer neural network. For the traditional NGs, we view both the VAE model and the VMF model as 3-layer neural networks. We apply K-FAC on these models to efficiently approximate the Fisher information matrices with respect to the model distributions, given the samples from the variational distributions.

To apply the K-FAC approximation, we will need to compute matrix multiplications and matrix inversions for some non-diagonal large square matrices (i.e. the $\bar{A}_{0,0}$ matrices in the K-FAC paper (Martens & Grosse, 2015) and some other matrices that are computed during the K-FAC approximation process). In order to make the algorithms faster, we applied low-rank approximations for some large matrices of these forms by sparse eigenvalue decompositions. All of these large matrices are positive semi-definite. For each such large matrix M , we keep only the $K \cdot \ln(\dim(M))$ dimensions of it with the largest eigenvalues and K is a hyperparameter that can be tuned.

For NG and VPNG, we applied the exponential moving averages for all matrices $\bar{A}_{i,i}$ and $\bar{G}_{i+1,i+1}$ (again, we use the notations in the K-FAC paper (Martens & Grosse, 2015)) to make the learning process more stable. We found that, by adding the exponential moving average technique, our VPNG performs similarly to the case without this technique, while the traditional natural gradient is much more stable and efficient. If we do not apply the exponential moving average technique, the traditional natural gradients will not perform well. As an example, we show the performances of all methods without the exponential moving average technique in the VAE experiment in Figure 5.

We found that NG and VPNG performed similarly with respect to the dampening factor μ , the exponential moving average decay parameter and the low-rank approximation function parameter K . However, different step sizes are needed to get the best performance from these two methods. We grid searched the step sizes and report the best one.