# Natural Conjugate Gradient in Variational Inference

Antti Honkela, Matti Tornio, Tapani Raiko, and Juha Karhunen

Adaptive Informatics Research Centre, Helsinki University of Technology

`http://www.cis.hut.fi/projects/bayes/`

E-mail: {antti.honkela, matti.tornio, tapani.raiko, juha.karhunen}@tkk.fi

**Summary:** Variational methods for approximate inference in machine learning often adapt a parametric probability distribution to optimize a given objective function. This view is especially useful when applying variational Bayes (VB) to models outside the conjugate-exponential family. For them, variational Bayesian expectation maximization (VB EM) algorithms are not easily available, and gradient-based methods are often used as alternatives. Traditional natural gradient methods use the Riemannian structure (or geometry) of the predictive distribution to speed up maximum likelihood estimation. We propose using the geometry of the variational approximating distribution instead to speed up a conjugate gradient method for variational learning and inference. The computational overhead is small due to the simplicity of the approximating distribution. Experiments with real-world speech data show significant speedups over alternative learning algorithms.

**Theory:** In previous machine learning algorithms based on natural gradients [1], the aim has been to use maximum likelihood to directly update the model parameters $\theta$ taking into account the geometry imposed by the predictive distribution for data $p(X|\theta)$. The resulting geometry is often very complicated as the effects of different parameters cannot be separated and the Fisher information matrix is relatively dense.

In this paper we propose using natural gradients for free energy minimisation in variational Bayesian learning using the simpler geometry of the approximating distributions $q(\theta|\xi)$. Because the approximations are often chosen to minimize dependencies between different parameters $\theta$, the resulting Fisher information matrix with respect to the variational parameters $\xi$ will be mostly diagonal and hence easy to invert.

While taking into account the structure of the approximation, plain natural gradient in this case ignores the structure of the model and the global geometry of the parameters $\theta$. This can be addressed by using conjugate gradients. Combining the natural gradient search direction with a conjugate gradient method yields our proposed *natural conjugate gradient (NCG)* method, which can also be seen as an approximation to the fully Riemannian conjugate gradient method.

**Experimental results:** The NCG algorithm was compared against conjugate gradient (CG) and natural gradient (NG) algorithms in learning a nonlinear state-space model [2]. The results for a number of datasets ranging from 200 to 500 samples of 21 dimensional speech spectrograms can be seen in Figure 1. The plain CG and NG methods were clearly slower than others and the maximum runtime of 24 hours was reached by most CG and some NG runs. NCG was clearly the fastest algorithm with the
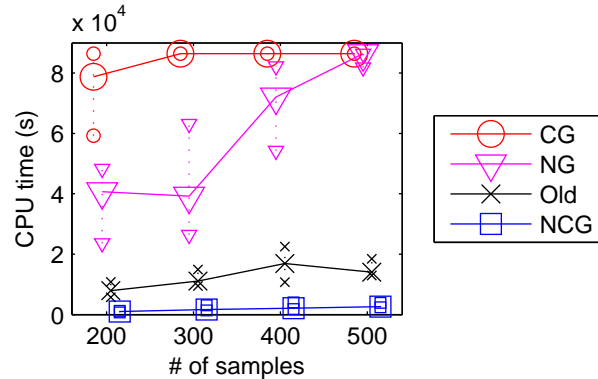


Fig. 1: Convergence speed of the natural conjugate gradient (NCG), the natural gradient (NG) and the conjugate gradient (CG) methods as well as the heuristic algorithm (Old) with different data sizes. The lines show median times with 25 % and 75 % quantiles shown by the smaller marks. The times were limited to at most 24 hours, which was reached by a number of simulations.

heuristic method of [2] between these extremes.

The results with a larger data set are very similar with NCG outperforming all alternatives by a factor of more than 10. For more details, please refer to the paper.

**Discussion:** The experiments in this paper show that the natural conjugate gradient method outperforms both conjugate gradient and natural gradient methods by a large margin. Considering univariate Gaussian distributions, the regular gradient is too strong for model variables with small posterior variance and too weak for variables with large posterior variance. The posterior variance of latent variables is often much larger than the posterior variance of model parameters and the natural gradient takes this into account in a very natural manner.

## References

[1] S. Amari. Natural gradient works efficiently in learning. *Neural Computation*, 10(2):251–276, 1998.

[2] H. Valpola and J. Karhunen. An unsupervised ensemble learning method for nonlinear dynamic state-space models. *Neural Computation*, 14(11):2647–2692, 2002.