

# The Variational Predictive Natural Gradient

Da Tang<sup>1</sup> Rajesh Ranganath<sup>2</sup>

<sup>1</sup>Columbia University

<sup>2</sup>New York University

June 12, 2019

- ▶ Latent variable models:  $p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})$ .

- ▶ Latent variable models:  $p(\mathbf{x}, \mathbf{z}; \theta) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z}; \theta)$ .
- ▶ Variational inference approximates the posterior through maximizing the ELBO:

$$\mathcal{L}(\lambda, \theta) = \mathbb{E}_q [\log p(\mathbf{x}|\mathbf{z}; \theta)] - \text{KL}(q(\mathbf{z}|\mathbf{x}; \lambda) || p(\mathbf{z})).$$

- ▶ Latent variable models:  $p(\mathbf{x}, \mathbf{z}; \theta) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z}; \theta)$ .
- ▶ Variational inference approximates the posterior through maximizing the ELBO:

$$\mathcal{L}(\boldsymbol{\lambda}, \theta) = \mathbb{E}_q [\log p(\mathbf{x}|\mathbf{z}; \theta)] - \text{KL}(q(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda})||p(\mathbf{z})).$$

- ▶  $q$ -Fisher Information  $F_q = \mathbb{E}_q [\nabla_{\boldsymbol{\lambda}} \log q(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda}) \cdot \nabla_{\boldsymbol{\lambda}} \log q(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda})^\top]$  (Hoffman et al., 2013) approximates the negative Hessian of the objective.
- ▶ The natural gradient:  $\nabla_{\boldsymbol{\lambda}}^{\text{NG}} \mathcal{L}(\boldsymbol{\lambda}) = F_q^{-1} \cdot \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda})$ .

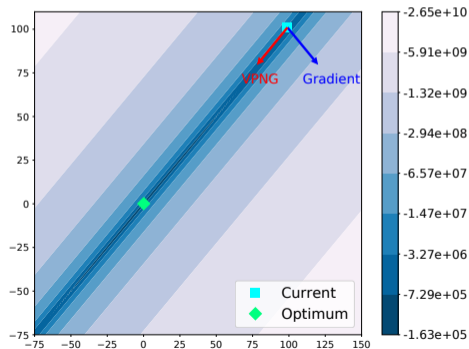
# Pathological Curvature of the ELBO

- ▶ The curvature of the ELBO may be pathological.

# Pathological Curvature of the ELBO

- ▶ The curvature of the ELBO may be pathological.
- ▶ Example: A bivariate Gaussian model with unknown mean and known covariance

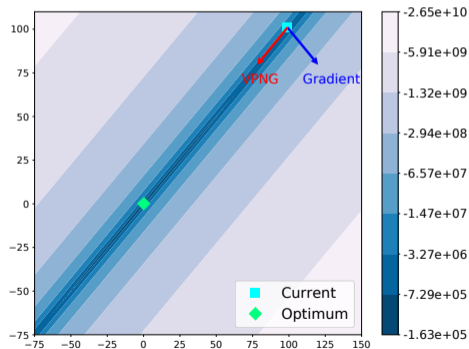
$$\Sigma = \begin{pmatrix} 1 & 1 - \varepsilon \\ 1 - \varepsilon & 1 \end{pmatrix}, 0 < \varepsilon \ll 1.$$



# Pathological Curvature of the ELBO

- ▶ The curvature of the ELBO may be pathological.
- ▶ Example: A bivariate Gaussian model with unknown mean and known covariance

$$\Sigma = \begin{pmatrix} 1 & 1 - \varepsilon \\ 1 - \varepsilon & 1 \end{pmatrix}, 0 < \varepsilon \ll 1.$$



- ▶ The natural gradient fails to help.

# The Natural Gradient is Insufficient

Limitations of the  $q$ -Fisher information:

- ▶ Approximates the Hessian of the objective well only when  $q(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda}) \approx p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})$ .
- ▶ Ignore the model likelihood  $p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})$  in computations.



# The Variational Predictive Fisher Information

- ▶ Construct a **positive definite** matrix that resembles the negative Hessian of the expected log-likelihood part  $\mathcal{L}^{\text{ll}} = \mathbb{E}_{q(\mathbf{z}|\mathbf{x};\lambda)} [\log p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})]$  of the ELBO.

# The Variational Predictive Fisher Information

- ▶ Construct a **positive definite** matrix that resembles the negative Hessian of the expected log-likelihood part  $\mathcal{L}^{\text{ll}} = \mathbb{E}_{q(\mathbf{z}|\mathbf{x};\lambda)} [\log p(\mathbf{x}|\mathbf{z};\theta)]$  of the ELBO.
- ▶ Reparameterize the variational distribution  $q$ :

$$\mathbf{z} = g(\mathbf{x}, \boldsymbol{\varepsilon}; \boldsymbol{\lambda}) \sim q(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda}) \iff \boldsymbol{\varepsilon} \sim s(\boldsymbol{\varepsilon}).$$

# The Variational Predictive Fisher Information

- ▶ Construct a **positive definite** matrix that resembles the negative Hessian of the expected log-likelihood part  $\mathcal{L}^{\text{ll}} = \mathbb{E}_{q(\mathbf{z}|\mathbf{x};\boldsymbol{\lambda})} [\log p(\mathbf{x}|\mathbf{z};\boldsymbol{\theta})]$  of the ELBO.
- ▶ Reparameterize the variational distribution  $q$ :

$$\mathbf{z} = g(\mathbf{x}, \boldsymbol{\varepsilon}; \boldsymbol{\lambda}) \sim q(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda}) \iff \boldsymbol{\varepsilon} \sim s(\boldsymbol{\varepsilon}).$$

- ▶ The variational predictive Fisher information:

$$F_r = \mathbb{E}_{\boldsymbol{\varepsilon}} [\mathbb{E}_{p(\mathbf{x}'|\mathbf{z}=g(\mathbf{x},\boldsymbol{\varepsilon};\boldsymbol{\lambda});\boldsymbol{\theta})} [\nabla_{\boldsymbol{\lambda},\boldsymbol{\theta}} \log p(\mathbf{x}'|\mathbf{z} = g(\mathbf{x}, \boldsymbol{\varepsilon}; \boldsymbol{\lambda}); \boldsymbol{\theta}) \cdot \nabla_{\boldsymbol{\lambda},\boldsymbol{\theta}} \log p(\mathbf{x}'|\mathbf{z} = g(\mathbf{x}, \boldsymbol{\varepsilon}; \boldsymbol{\lambda}); \boldsymbol{\theta})^\top]],$$

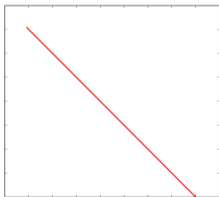
exactly the “expected” Fisher information of the *reparameterized predictive distribution*  $p(\mathbf{x}'|\mathbf{z} = g(\mathbf{x}, \boldsymbol{\varepsilon}; \boldsymbol{\lambda}); \boldsymbol{\theta})$ .

# The Variational Predictive Fisher Information

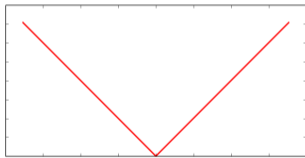
- ▶ Variational predictive Fisher captures the curvature of variational inference.

# The Variational Predictive Fisher Information

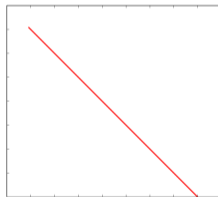
- ▶ Variational predictive Fisher captures the curvature of variational inference.
- ▶ Matrix spectrum comparison (for the bivariate Gaussian example):



(d) Precision mat  $\Sigma^{-1}$



(e)  $q$ -Fisher info  $F_q$



(f) Our Fisher info  $F_r$

# The Variational Predictive Natural Gradient

- ▶ The variational predictive natural gradient (VPNG):

$$\nabla_{\lambda, \theta}^{\text{VPNG}} \mathcal{L} = F_r^{-1} \cdot \nabla_{\lambda, \theta} \mathcal{L}(\lambda, \theta).$$

# The Variational Predictive Natural Gradient

- ▶ The variational predictive natural gradient (VPNG):

$$\nabla_{\lambda, \theta}^{\text{VPNG}} \mathcal{L} = F_r^{-1} \cdot \nabla_{\lambda, \theta} \mathcal{L}(\lambda, \theta).$$

- ▶ In practice, use Monte Carlo estimations to approximate  $F_r$  and add a small dampening parameter to ensure invertibility.

# Experiments: Bayesian Logistic Regression

- ▶ Tested on synthetic data with high correlations.
- ▶ Empirical results:

| Method   | Train AUC                           | Test AUC                            |
|----------|-------------------------------------|-------------------------------------|
| Gradient | $0.734 \pm 0.017$                   | $0.718 \pm 0.022$                   |
| NG       | $0.744 \pm 0.043$                   | $0.751 \pm 0.047$                   |
| VPNG     | <b><math>0.972 \pm 0.011</math></b> | <b><math>0.967 \pm 0.011</math></b> |

Table: Bayesian Logistic regression AUC



# Experiments: VAE and VMF

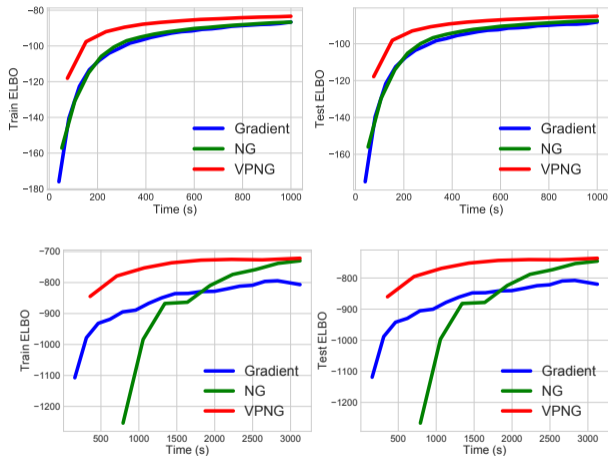


Figure: Learning curves of variational autoencoders (upper) and variational matrix factorization (lower) on real datasets.

## Conclusion and Future Work

- ▶ The VPNG corrects for curvature in the objective between the parameters in variational inference.

## Conclusion and Future Work

- ▶ The VPNG corrects for curvature in the objective between the parameters in variational inference.
- ▶ Future work includes extending to general Bayesian networks with multiple stochastic layers.

# Thanks!

Poster #234

Code available at <https://github.com/datang1992/VPNG>.