

## Introduction

### Issues:

- Existing variational inference (VI) methods, e.g., black-box VI (BBVI) [2], require significant changes to existing deep-learning implementation and demand large memory.

### Contributions:

- For Gaussian variational approximations, we propose a method called **Vprop** which can be implemented with a slight modification to **RMSprop** code.
- Vprop simplifies the implementation of **BBVI** and also reduces its memory requirement by half.
- Vprop is a **natural-gradient** method for VI and also related to **Newton's method**.

### Algorithm 1 Vprop-1 for VI

- $\theta \leftarrow \mu + \epsilon \cdot \sqrt{\mathbf{s} + \lambda}$ ,  $\epsilon \sim \mathcal{N}(\epsilon|0, \mathbf{I})$
- $\mathbf{g} \leftarrow \nabla_{\theta} f(\theta)$
- $\mathbf{s} \leftarrow (1 - \beta)\mathbf{s} + \beta(\mathbf{g} * \mathbf{g})$
- $\mu \leftarrow \mu - \beta((\mathbf{g} + \lambda\mu) / (\mathbf{s} + \lambda))$

(a) Vprop-1 minimizes  $\mathcal{L}(\mu, \sigma) = \mathbb{E}_q[f(\theta) + \log p(\theta)/q(\theta)]$ .

### Algorithm 2 RMSprop for MLE

- $\theta \leftarrow \mu$
- $\mathbf{g} \leftarrow \nabla_{\theta} f(\theta)$
- $\mathbf{s} \leftarrow (1 - \beta)\mathbf{s} + \beta(\mathbf{g} * \mathbf{g})$
- $\mu \leftarrow \mu - \alpha(\mathbf{g} / \sqrt{\mathbf{s} + \delta})$

(b) RMSprop minimizes  $f(\theta) := -\log p(\mathbf{y}|\mathbf{X}, \theta)$ .

## Optimization Algorithms & VI

### Maximum Likelihood Estimation (MLE):

We can perform **MLE** by minimizing  $f(\theta) = -\log p(\mathbf{y}|\mathbf{X}, \theta)$  using RMSprop [1], as shown in Algorithm 2:

$$\text{RMSprop: } \begin{aligned} \theta_{t+1} &= \theta_t - \alpha_t (\mathbf{s}_{t+1} + \delta \mathbf{1})^{-\frac{1}{2}} \circ \nabla_{\theta} f(\theta_t), \\ \mathbf{s}_{t+1} &= (1 - \beta_t)\mathbf{s}_t + \beta_t (\nabla_{\theta} f(\theta_t))^2, \end{aligned}$$

where  $\mathbf{s}$  is the adaptive step-size scaling vector.

### Gaussian Variational Inference:

In contrast, in **Gaussian VI**, we approximate the posterior distribution with a Gaussian distribution  $q(\theta) := \mathcal{N}(\theta|\mu, \text{diag}(\sigma^2))$  by maximizing a variational lower-bound:

$$\log p(\mathbf{y}|\mathbf{X}) = \max_{\mu, \sigma} \mathbb{E}_q[\log p(\mathbf{y}|\mathbf{X}, \theta) + \log p(\theta) - \log q(\theta)] := \mathcal{L}(\mu, \sigma)$$

where  $p(\theta) = \mathcal{N}(\theta|0, \mathbf{I}/\lambda)$  is the Gaussian prior.

### Black-box Variational Inference:

**Black-box VI** (BBVI) [2] optimizes  $\mu$  and  $\sigma$  by the following stochastic gradient-descent algorithm (or adaptive-gradient variants of it):

$$\text{BBVI: } \begin{aligned} \mu_{t+1} &= \mu_t + \eta_t \left[ \widehat{\nabla}_{\mu} \mathcal{L}_t \right], \\ \sigma_{t+1} &= \sigma_t + \rho_t \left[ \widehat{\nabla}_{\sigma} \mathcal{L}_t \right]. \end{aligned}$$

However, BBVI update has two issues.

- Memory:** The number of parameters is doubled for adaptive-gradient since it needs to store scaling vectors for  $\mu$  and  $\sigma$ .
- Implementation:** Computing the gradients w.r.t.  $\sigma$  requires a different implementation from that of existing deep learning.

### Conjugate-computation Variational Inference:

We propose Vprop to solve these issues by using a natural gradient method called **conjugate-computation VI** (CVI) [3]:

$$\text{CVI: } \begin{aligned} \mu_{t+1} &= \mu_t + \beta_t \sigma_{t+1}^2 \circ [\nabla_{\mu} \mathcal{L}_t], \\ \sigma_{t+1}^{-2} &= \sigma_t^{-2} - 2\beta_t [\nabla_{\sigma^2} \mathcal{L}_t]. \end{aligned}$$

## Vprop

### Derivation of Vprop:

Vprop is derived from the CVI update in two steps.

- Apply Bonnet's and Price's theorem:

$$\begin{aligned} \nabla_{\mu} \mathcal{L} &= -\mathbb{E}_q[\nabla_{\theta} f(\theta)] - \lambda\mu, \\ \nabla_{\sigma} \mathcal{L} &= -\frac{1}{2} \mathbb{E}_q[\text{diag}(\nabla_{\theta}^2 f(\theta))] - \frac{1}{2}\lambda\mathbf{1} + \frac{1}{2}\sigma^{-2}. \end{aligned}$$

The Hessian is approximated by a **Gauss-Newton approximation**.

By defining  $\mathbf{s}_t := \sigma_t^{-2} - \lambda\mathbf{1}$  and using the above two steps in the CVI update, we obtain Vprop:

$$\text{Vprop: } \begin{aligned} \mu_{t+1} &= \mu_t - \beta_t (\mathbf{s}_{t+1} + \lambda\mathbf{1})^{-1} \circ \{\mathbb{E}_{q_t}[\nabla_{\theta} f(\theta)] + \lambda\mu_t\}, \\ \mathbf{s}_{t+1} &= (1 - \beta_t)\mathbf{s}_t + \beta_t \mathbb{E}_{q_t}[(\nabla_{\theta} f(\theta))^2]. \end{aligned}$$

The expectations can be approximated by using one Monte-Carlo (MC) sample giving us Vprop-1:

$$\text{Vprop-1: } \begin{aligned} \theta_t &\sim \mathcal{N}(\theta|\mu_t, \text{diag}(\mathbf{s}_t^{-2})), \\ \mu_{t+1} &= \mu_t - \beta_t (\mathbf{s}_{t+1} + \lambda\mathbf{1})^{-1} \circ \{\nabla_{\theta} f(\theta_t) + \lambda\mu_t\}, \\ \mathbf{s}_{t+1} &= (1 - \beta_t)\mathbf{s}_t + \beta_t (\nabla_{\theta} f(\theta_t))^2. \end{aligned}$$

### Comparison between Vprop-1 and RMSprop:

Vprop-1 and RMSprop are similar but have 3 significant differences:

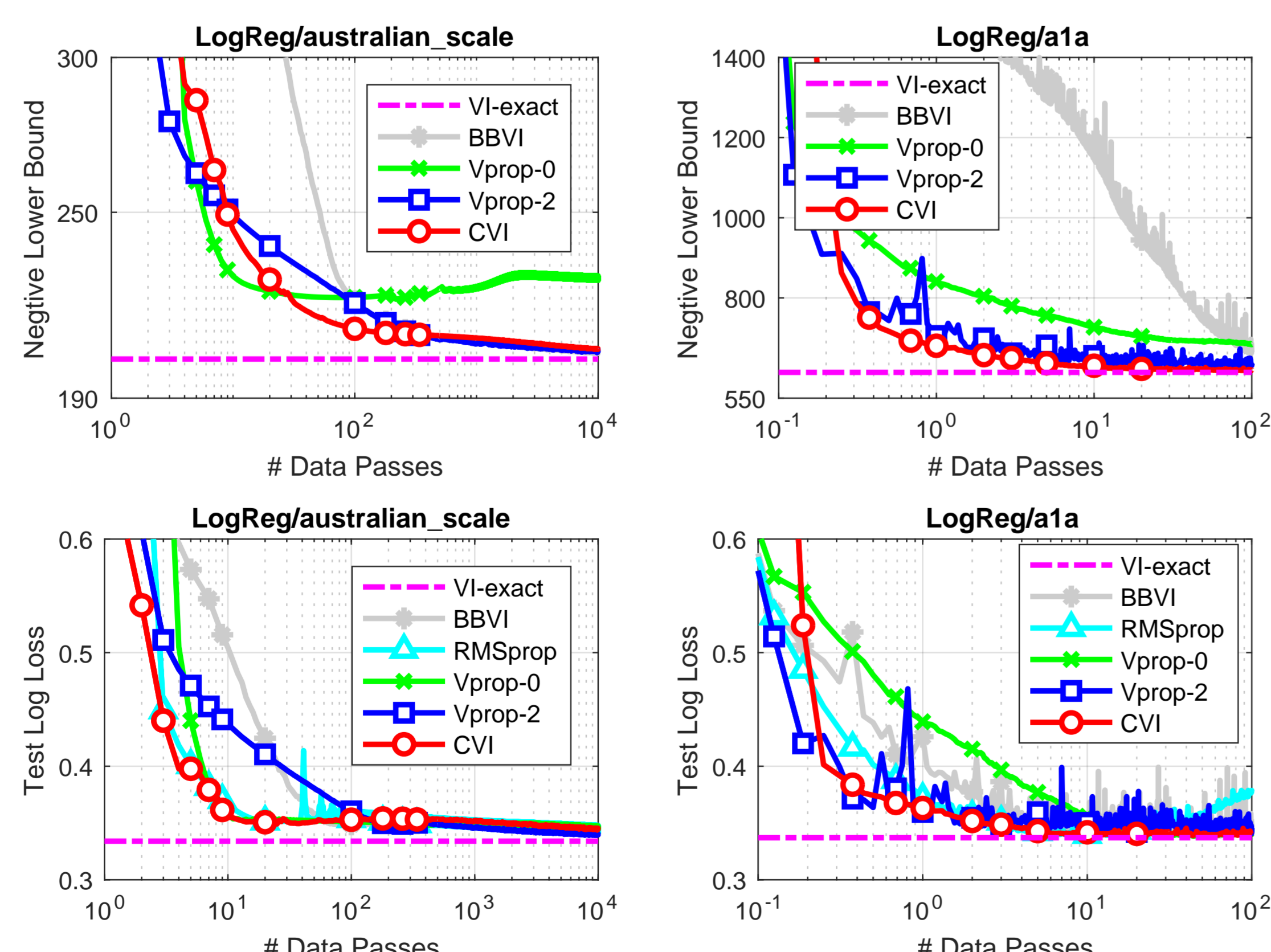
- Vprop-1 computes the gradients at  $\theta$  sampled from  $\mathcal{N}(\theta|\mu_t, \text{diag}(\mathbf{s}_t^{-2}))$  but RMSprop computes gradients at  $\theta = \mu_t$ .
- RMSprop raise the scaling vector to a power of  $\frac{1}{2}$ , while Vprop-1 does not.
- Vprop-1 adds  $\lambda\mu$  to the gradient of  $\mu$  update.

### Connections to natural gradient and Newton's method:

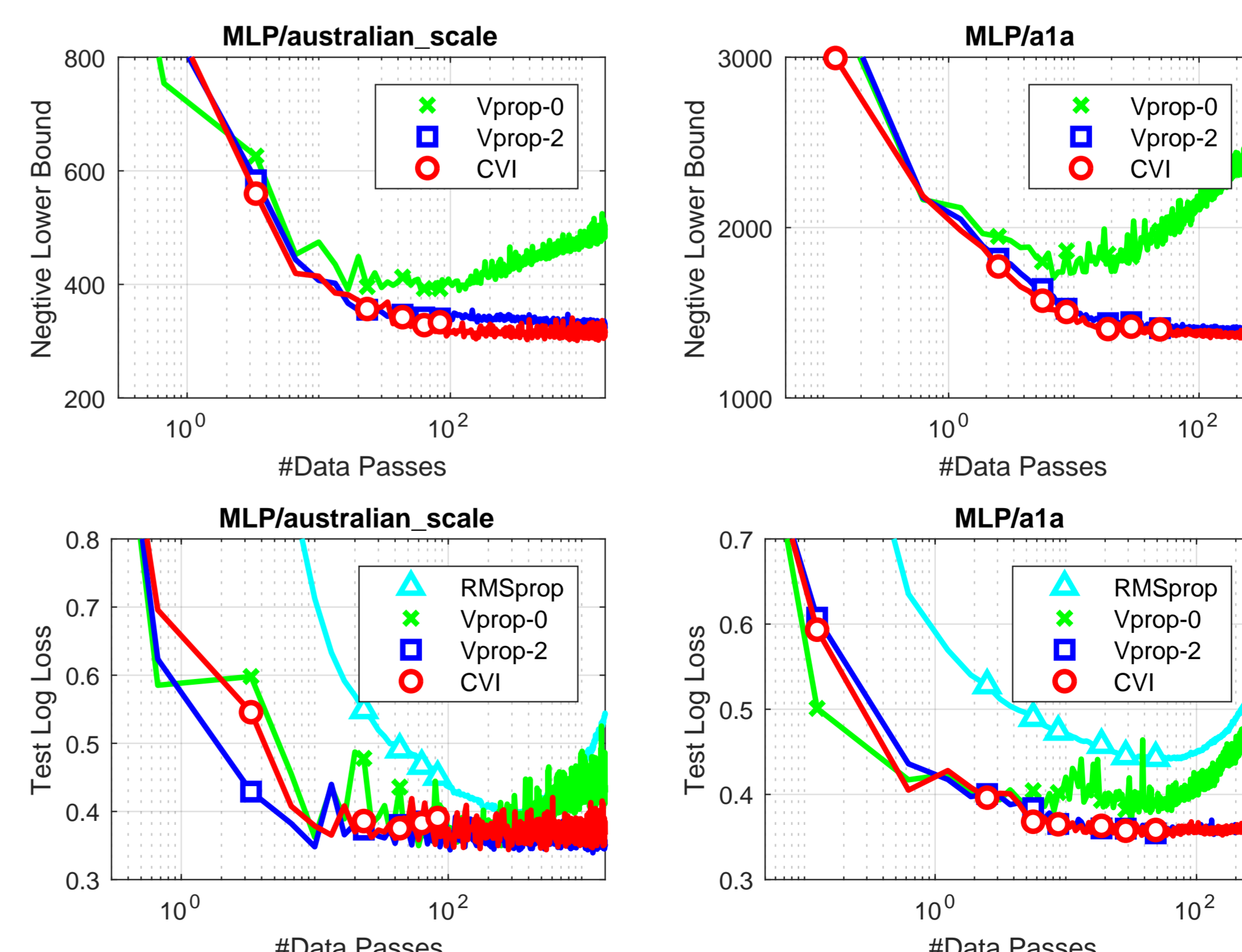
For full Gaussian  $q(\theta)$ , using similar derivations result in a variant of **Newton's method**. The resulting update also resembles an **online natural gradient** descent proposed by Ollivier 2017 [4].

## Experiments

### Logistic regression



### Multilayer-perceptron



## References

- [1] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural Networks for Machine Learning: Lecture 6a.
- [2] Rajesh Ranganath, Sean Gerrish, and David M Blei. Black box variational inference. *AISTATS*, 2014.
- [3] Mohammad Emtiyaz Khan and Wu Lin. Conjugate-computation variational inference. *AISTATS*, 2017.
- [4] Yann Ollivier. Online Natural Gradient as a Kalman Filter. *arXiv*, 2017.