

Accelerating Natural Gradient with Higher-Order Invariance

Yang Song, Jiaming Song, Stefano Ermon
Stanford University
{yangsong, tsong, ermon}@cs.stanford.edu

ABSTRACT

The finite step size in natural gradient descent makes the optimization trajectory not invariant to model reparameterizations. We propose several ways to improve its invariance.

- We propose to measure the invariance of optimization methods by comparing their convergence to idealized invariant trajectories.
- We propose to use midpoint integrators for improving natural gradient optimization
- We introduce geodesic corrected updates, including a faster version which has comparable time complexity to vanilla natural gradient. We prove the convergence for both types of geodesic-corrected updates.
- Experiments demonstrate faster convergence of our proposed algorithms in supervised learning and reinforcement learning applications.

Paper:
<https://arxiv.org/abs/1803.01273>

Code:
https://github.com/ermongroup/higher_order_invariance

INVARIANCE AND DIFFERENTIAL GEOMETRY

Einstein's summation convention: A repeated pair of index variables—one as superscript and one as subscript—indicates summation over it.

$$a^\mu b_\mu \stackrel{\text{def}}{=} \sum_{\mu=1}^n a^\mu b_\mu$$

Manifold: A smooth space M where at each point $p \in M$ you can define tangent spaces T_p and cotangent spaces T_p^* , both are Euclidean spaces.

Coordinates: Real numbers used to describe a point on the manifold, or vectors and tensors on tangent/cotangent spaces. Coordinates are w.r.t. a coordinate system. The same entity has different coordinates in different coordinate systems.

Examples of coordinates:

A point: $p = (\theta^1, \theta^2, \dots, \theta^n) \in M$

Bases of T_p : $\partial/\partial\theta^1, \partial/\partial\theta^2, \dots, \partial/\partial\theta^n$, abbr $\partial_\mu = \partial/\partial\theta^\mu$

A vector in T_p : $\mathbf{a} = (a^1, a^2, \dots, a^n) = a^\mu \partial_\mu \in T_p$

Bases of T_p^* : $d\theta^1, d\theta^2, \dots, d\theta^n$

A covector in T_p^* : $\mathbf{a}^* = (a_1, a_2, \dots, a_n) = a_\mu d\theta^\mu \in T_p^*$

A tensor in $T_p^* \otimes T_p^*$: $\mathbf{g} = \begin{pmatrix} g_{11} & \dots & g_{1n} \\ \vdots & \ddots & \vdots \\ g_{n1} & \dots & g_{nn} \end{pmatrix} = g_{\mu\nu} d\theta^\mu \otimes d\theta^\nu$

Coordinate transformations: How should coordinates transform when the coordinate system changes from $(\theta^1, \theta^2, \dots, \theta^n)$ to $(\xi^1, \xi^2, \dots, \xi^n)$? Depending on whether it is a **superscript** or **subscript**, each rank should change as

$$a^{\mu'} = a^\mu \frac{\partial \xi^{\mu'}}{\partial \theta^\mu}$$

$$a_{\mu'} = a_\mu \frac{\partial \theta^\mu}{\partial \xi^{\mu'}}$$

Geodesic equation: The equation governing coordinates of points in a geodesic line.

$$\frac{d^2 \gamma^\mu}{dt^2} + \Gamma_{\alpha\beta}^\mu \frac{d\gamma^\alpha}{dt} \frac{d\gamma^\beta}{dt} = 0$$

$$\Gamma_{\alpha\beta}^\mu \stackrel{\text{def}}{=} \frac{1}{2} g^{\mu\nu} (\partial_\alpha g_{\nu\beta} + \partial_\beta g_{\nu\alpha} - \partial_\nu g_{\alpha\beta})$$

Exponential map: If we follow the curve $\gamma(t)$ from $p = \gamma(0)$ with initial direction $v = \frac{d\gamma(t=0)}{dt}$ for a unit time $\Delta t = 1$, which point can we reach?

$$\text{Exp}(p, v) \stackrel{\text{def}}{=} \gamma(1)$$

Re-scaling gives us $\text{Exp}(p, hv) = \gamma(h)$.

Invariance: An equation written in coordinates should not change due to coordinate transformations.

REVISITING NATURAL GRADIENT

Model family as a manifold:

$$r_\theta(x, t) = p_\theta(t|x)q(x) \quad \rightarrow \text{Point on a manifold}$$

$$\theta \quad \rightarrow \text{Coordinates}$$

Natural gradient ODE:

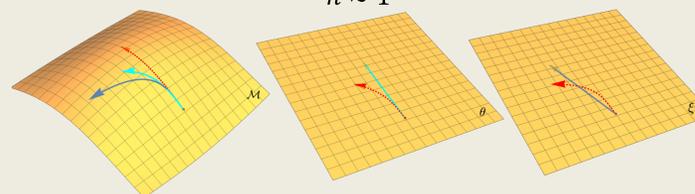
$$\text{Superscript} \rightarrow \frac{d\theta^\mu}{dt} = -\lambda g^{\mu\nu} \partial_\nu L(r_\theta), \quad \leftarrow \text{Superscript}$$

where $L(r_\theta)$ is the loss function, $g^{\mu\nu}$ is the inverse of Fisher information metric.

Naïve natural gradient update:

$$\theta_{k+1}^\mu \leftarrow \theta_k^\mu - h \lambda g^{\mu\nu} \partial_\nu L(r_{\theta_k})$$

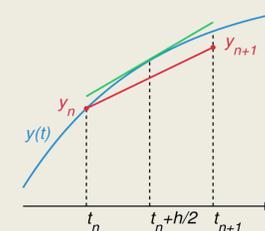
$$h \approx 1$$



HIGHER-ORDER INTEGRATORS

Motivations: Since the natural gradient ODE is invariant, more accurate integrators can give more invariant optimization trajectories.

Midpoint integrator:



$$\theta_{k+1/2}^\mu \leftarrow \theta_k^\mu - \frac{1}{2} h \lambda g^{\mu\nu} \partial_\nu L(r_{\theta_k})$$

$$\theta_{k+1}^\mu \leftarrow \theta_k^\mu - h \lambda g^{\mu\nu} \partial_\nu L(r_{\theta_{k+1/2}})$$

GEODESIC CORRECTIONS

Motivations: Geodesics are invariant.

Riemannian Euler Method: Exactly invariant solver, even with finite step sizes.

$$\theta_{k+1}^\mu \leftarrow \text{Exp}(\theta_k^\mu, -h \lambda g^{\mu\nu} \partial_\nu L(r_{\theta_k}))$$

Geodesic corrections: Use the geodesic equation to approximately compute exponential map.

$$\theta_{k+1}^\mu \leftarrow \theta_k^\mu + h \frac{d\gamma_k^\mu(t=0)}{dt} - \frac{1}{2} h^2 \Gamma_{\alpha\beta}^\mu \frac{d\gamma_k^\alpha(t=0)}{dt} \frac{d\gamma_k^\beta(t=0)}{dt}$$

$$\frac{d\gamma_k^\mu(t=0)}{dt} \equiv -\lambda g^{\mu\nu} \partial_\nu L(r_{\theta_k})$$

Computations: Tensor products can be done by forward-mode autodifferentiation. $g^{\mu\nu}$ can be obtained by approximations (e.g. truncated CG, KFAC).

Faster geodesic corrections: Approximate

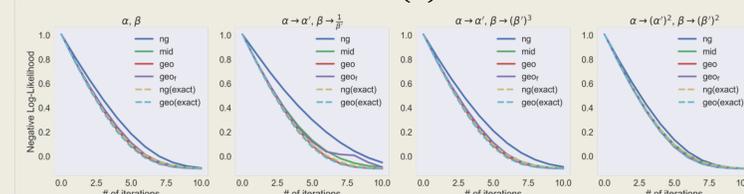
$$\frac{d\gamma_k^\mu(t=0)}{dt} \text{ in } \Gamma_{\alpha\beta}^\mu \frac{d\gamma_k^\alpha(t=0)}{dt} \frac{d\gamma_k^\beta(t=0)}{dt} \text{ with } \frac{(\theta_k^\mu - \theta_{k-1}^\mu)}{h}$$

Theorem 1 (informal): As $h \rightarrow 0$, the Euler integrator used in naive natural gradient update converges to Riemannian Euler method's solution in 1st order, while both kinds of geodesic corrected integrators converge in 2nd order.

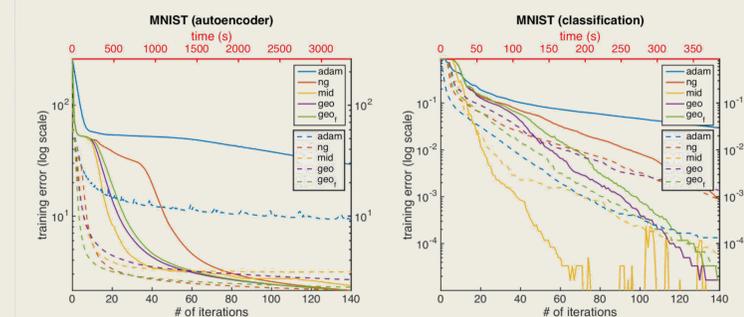
EXPERIMENTS

Invariance: Fitting a univariate Gamma distribution with different parameterizations.

$$p(x|\alpha, \beta) \stackrel{\text{def}}{=} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$$



Training Deep Neural Networks:



Policy Gradient Reinforcement Learning:

