

Gradient Descent

Yutao Chen

✂ Oct 03 2025

✂ Oct 08 2025

Contents

Lagrangian Method	1
Riemannian Manifold	2

In machine learning, we are often tasked with updating the model parameters $\theta \in \mathbb{R}^d$ to minimize some objective function $J(\theta) : \mathbb{R}^d \mapsto \mathbb{R}$.

Gradient descent tells us to update θ in the direction of the negative gradient, that is

$$\theta_{t+1} = \theta_t - \eta \nabla J(\theta_t),$$

where η is commonly known as the *learning rate*.

Lagrangian Method

Directly minimizing the black-box objective $J(\theta)$ is intractable in general. We instead consider a linear approximation of $J(\theta)$ around θ_t using the first-order Taylor expansion:

$$J(\theta) \approx J(\theta_t) + \nabla J(\theta_t)^\top (\theta - \theta_t).$$

We minimize the linear approximation of $J(\theta)$ with a squared Euclidean distance penalty $\|\theta - \theta_t\|_2^2$:

$$\arg \min_{\theta} J(\theta_t) + \nabla J(\theta_t)^\top (\theta - \theta_t) + (2\eta)^{-1} \|\theta - \theta_t\|_2^2,$$

where the learning rate η controls the *inverse* strength of the penalty. Intuitively, the penalty discourages large updates as the linear approximation might only hold in the vicinity of θ_t .

Setting the gradient of the optimization objective above to zero, we have

$$\begin{aligned} 0 &= \nabla J(\theta_t) + \eta^{-1}(\theta^* - \theta_t) \\ \implies \theta^* &= \theta_t - \eta \nabla J(\theta_t). \end{aligned}$$

Riemannian Manifold

By far, we have assumed a Euclidean parameter space Θ such that $\forall \theta, \theta' \in \Theta$, the squared distance between two parameters $D(\theta, \theta')$ is

$$D(\theta, \theta') = \|\theta - \theta'\|_2^2 = (\theta - \theta')^\top (\theta - \theta').$$

However, in general the parameter space Θ is not Euclidean. For example, one can parametrize the variance of a univariate Gaussian distribution by letting $\theta = \sigma^2$, $\theta = \sigma$, or $\theta = \log \sigma$ (and more). Depending the chosen parametrization, updating θ by the same constant c

$$\theta_{t+1} = \theta_t - c$$

will result in very different distributions.

More generally, one can consider the parameter space as a *Riemannian manifold*. The squared distance between two parameters θ and θ' on a Riemannian manifold¹ is

$$D(\theta, \theta') = (\theta - \theta')^\top G(\theta)(\theta - \theta'),$$

where the *Riemannian metric tensor* $G(\theta) \in \mathbb{R}^{d \times d}$ is a symmetric and positive definite matrix depending on θ . Note that the Euclidean space corresponds to a special case where $G(\theta) = I$.

Amari (1998) showed that the steepest descent on a Riemannian manifold is given by

$$\theta_{t+1} = \theta_t - \eta G^{-1}(\theta_t) \nabla J(\theta_t).$$

Particularly when $G(\theta)$ is the Fisher information matrix (FIM), this is known as the [[Natural Gradient Descent](#)].

REFERENCES

Amari, S.-I. (1998). Natural gradient works efficiently in learning. *Neural Computation*, 10(2), 251–276. <https://doi.org/10.1162/089976698300017746>

¹Strictly speaking, in the tangent space of the Riemannian manifold at θ .