

Numerical Optimization 06: 1st order methods

Qiang Zhu

University of Nevada Las Vegas

July 14, 2020

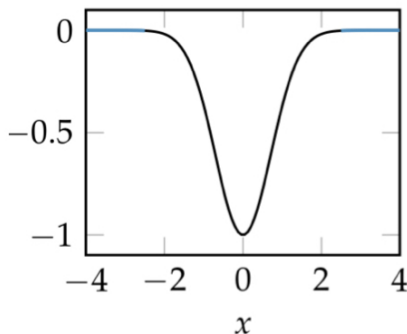
Overview

- 1 Gradient methods with fixed learning rate
- 2 Momentum
- 3 Nesterov Momentum
- 4 Adagrad
- 5 RMSProp and Adadelta
- 6 Adam
- 7 Summary

Gradient methods with fixed learning rate

We have talked about steepest descent and conjugate gradient methods, which usually work with the line search methods. Alternatively, it is popular to use the fixed learning rate method based on Gradient descent. However, the standard version will take a long time to traverse a nearly flat surface. Several methods have been proposed. They are commonly used in machine learning neural networks training.

- Momentum
- Nesterov Momentum
- Adagrad
- RMSProp
- Adadelata
- Adam



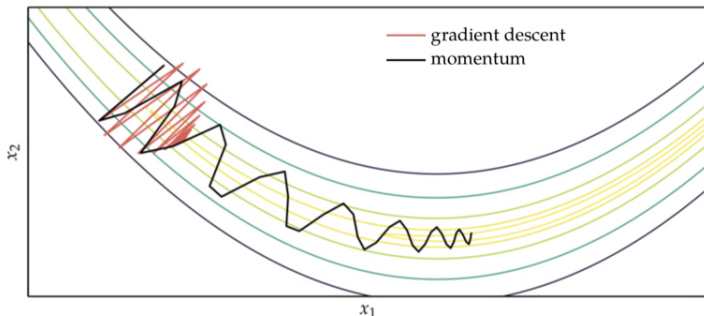
Momentum

Allowing momentum to accumulate is one way to speed the progress. Thus we can modify the gradient descent to incorporate momentum.

$$\mathbf{v}^{k+1} = \beta \mathbf{v}^k + \alpha \mathbf{g}^k$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \mathbf{v}^{k+1}$$

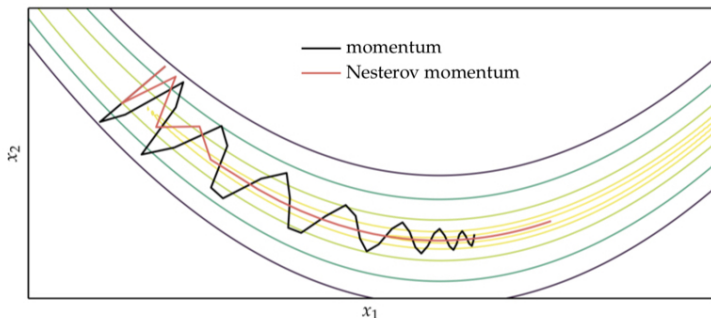
When $\beta=0$, it is gradient descent.



Nesterov Momentum

One issue of momentum is that the steps do not slow down enough at the bottom of a valley and it tends to **overshoot the valley**. Nesterov Momentum remedies the issue by the following updates,

$$\begin{aligned}\mathbf{v}^{k+1} &= \beta \mathbf{v}^k + \alpha \nabla f(\mathbf{x} + \beta \mathbf{v}^k) \\ \mathbf{x}^{k+1} &= \mathbf{x}^k + \mathbf{v}^{k+1}\end{aligned}$$



Adagrad

Momentum and Nesterov Momentum update all components of x with the same learning rate. The adaptive subgradient method (Adagrad), adapts a learning rate for each one in x .

$$x_i^{k+1} = x_i^k - \frac{\alpha}{\epsilon + \sqrt{s_i^k}} g^k$$
$$s_i^k = \sum_{j=1}^k \left(g_i^j \right)^2$$

where ϵ is a small value on the order of $1e-8$, to prevent the case of division by zero. Adagrad is far less sensitive to the learning rate α .

RMSProp and Adadelta

In Adagrad, the learning rate may monotonically decrease. To prevent this, **RMSprop** maintains a decaying average of squared gradients.

$$\mathbf{s}^{k+1} = \gamma \mathbf{s}^k + (1 - \gamma)(\mathbf{g}^k \odot \mathbf{g}^k)$$

where γ is between 0 and 1, and usually is 0.9.

$$\begin{aligned} x_i^{k+1} &= x_i^k - \frac{\alpha}{\epsilon + \sqrt{s_i^k}} g_i^k \\ &= x_i^k - \frac{\alpha}{\epsilon + \text{RMS}(g_i)} g_i^k \end{aligned}$$

While in **Adadelta**, an exponentially decaying average is used,

$$x_i^{k+1} = x_i^k - \frac{\text{RMS}(\delta x_i)}{\epsilon + \text{RMS}(g_i)} g_i^k$$

Adam

The adaptive moment estimation (Adam) is so far the most widely used optimization method in neural network training. It stores both an exponentially decaying squared gradient like RMSProp and Adadelta, but also an exponentially decaying gradient like momentum.

$$\mathbf{v}^{k+1} = \gamma_v \mathbf{v}^k + (1 - \gamma_v) \mathbf{g}^k$$

$$\mathbf{s}^{k+1} = \gamma_s \mathbf{s}^k + (1 - \gamma_s) \left(\mathbf{g}^k \odot \mathbf{g}^k \right)$$

$$\hat{\mathbf{v}}^{k+1} = \mathbf{v}^{k+1} / (1 - \gamma_v^k)$$

$$\hat{\mathbf{s}}^{k+1} = \mathbf{s}^{k+1} / (1 - \gamma_s^k)$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha \hat{\mathbf{v}}^{k+1} / \left(\epsilon + \sqrt{\hat{\mathbf{s}}^{k+1}} \right)$$

Summary

- Descent methods with momentum build up progress in favorable directions
- A wide variety of accelerated descent methods use special techniques to speed up descent