

On the momentum term in gradient descent learning algorithms

Daniel Richards Ravi Arputharaj , Adhithyan Kalaivanan

KTH

Formulation

A cornerstone of deep learning is non-convex optimization, and gradient descent is often the preferred method. Conventionally, a "momentum" term is added in the parameter update, to improve convergence rate and prevent being trapped in local minima.

$$\Delta w_t = -\epsilon \nabla_w \mathbf{L}(w) + p \Delta w_{t-1} \quad (1)$$

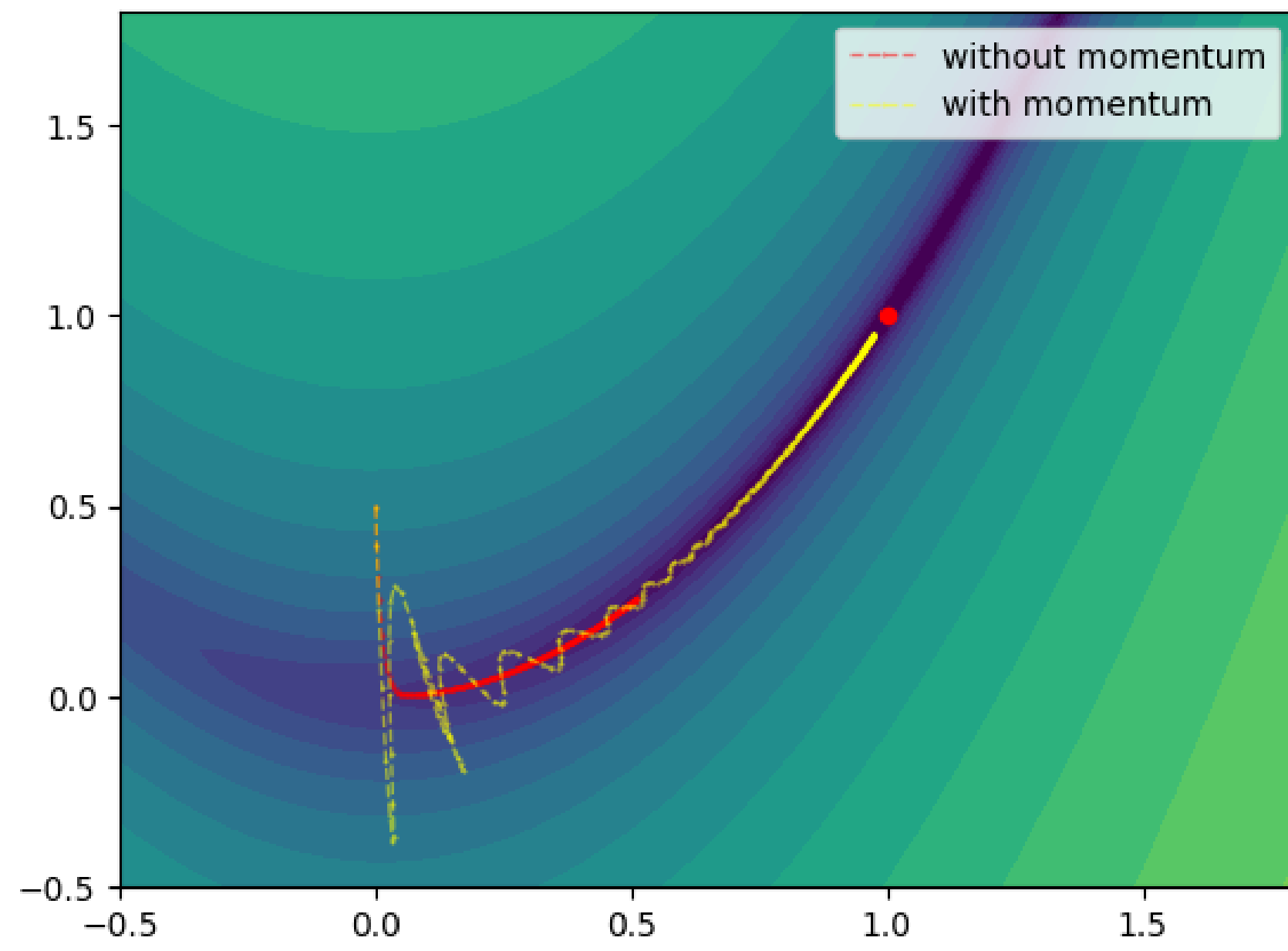


Figure 1. Effect of using momentum

An intuitive visual analogy is a ball rolling down the landscape of loss. Its inertia serves as an accelerator, helping speed of convergence.

Physical Analogy

Consider a point mass with mass m in a conservative force field with potential $\mathcal{E}(\mathbf{w})$ and friction coefficient μ [1]. The system is described by,

$$m \frac{d^2 \mathbf{w}}{dt^2} + \mu \frac{d\mathbf{w}}{dt} = -\nabla_w \mathcal{E}(\mathbf{w}) \quad (2)$$

where \mathbf{w} is the coordinate vector of the particle. Discretizing and re-arranging,

$$\mathbf{w}_{t+\Delta t} - \mathbf{w}_t = -\frac{(\Delta t)^2}{m + \mu \Delta t} \nabla_w \mathcal{E}(\mathbf{w}) + \frac{m}{m + \mu \Delta t} (w_t - w_{t-\Delta t}) \quad (3)$$

Comparing with equation 1, we see momentum term does act as inertia

$$\epsilon = \frac{(\Delta t)^2}{m + \mu \Delta t} \quad (4)$$

$$p = \frac{m}{m + \mu \Delta t} \quad (5)$$

Close to minima $_0$, equation 2 is approximated by,

$$m \frac{d^2 w}{dt^2} + \mu \frac{dw}{dt} = -\mathbf{H}(\mathbf{w} - \mathbf{w}_0) \quad (6)$$

where \mathbf{H} is the Hessian matrix. As it is p.s.d. $\mathbf{H} = QKQ^T$, with $K = \text{diag}(k_i | i = 1, 2, 3, \dots, n)$, where n is number of weights. Setting $\mathbf{w}' = Q^T \mathbf{w}$ then,

$$m \frac{d^2 w'_i}{dt^2} + \mu \frac{dw'_i}{dt} = -k_i w'_i \quad (7)$$

which resembles a set of uncoupled damped harmonic oscillators.

Analysis

The general solution to equation 7 is given by

$$w'_i(t) = c_1 e^{\lambda_{i,1} t} + c_2 e^{\lambda_{i,2} t} \quad (8)$$

$$\lambda_{i,1,2} = -\frac{\mu}{2m} \pm \sqrt{\frac{\mu}{m} \left(\frac{\mu}{4m} - \frac{k_i}{\mu} \right)} \quad (9)$$

If there were no momentum, the solution is simply

$$w'_i(t) = c_0 e^{\lambda_{i,0} t}, \lambda_{i,0} = -k_i / \mu \quad (9)$$

Observations:

- Convergence guaranteed: $\text{Re}(\lambda_{i,2}) \leq \text{Re}(\lambda_{i,1}) < 0$
- Momentum speed-up: $|\text{Re}(\lambda_{i,1})| > \text{Re}(|\lambda_{i,0}|)$, if and only if $k_i < \mu^2/2m$
- Setting $|\text{Re}(\lambda_{i,1})| := \alpha |\text{Re}(\lambda_{i,0})|$,

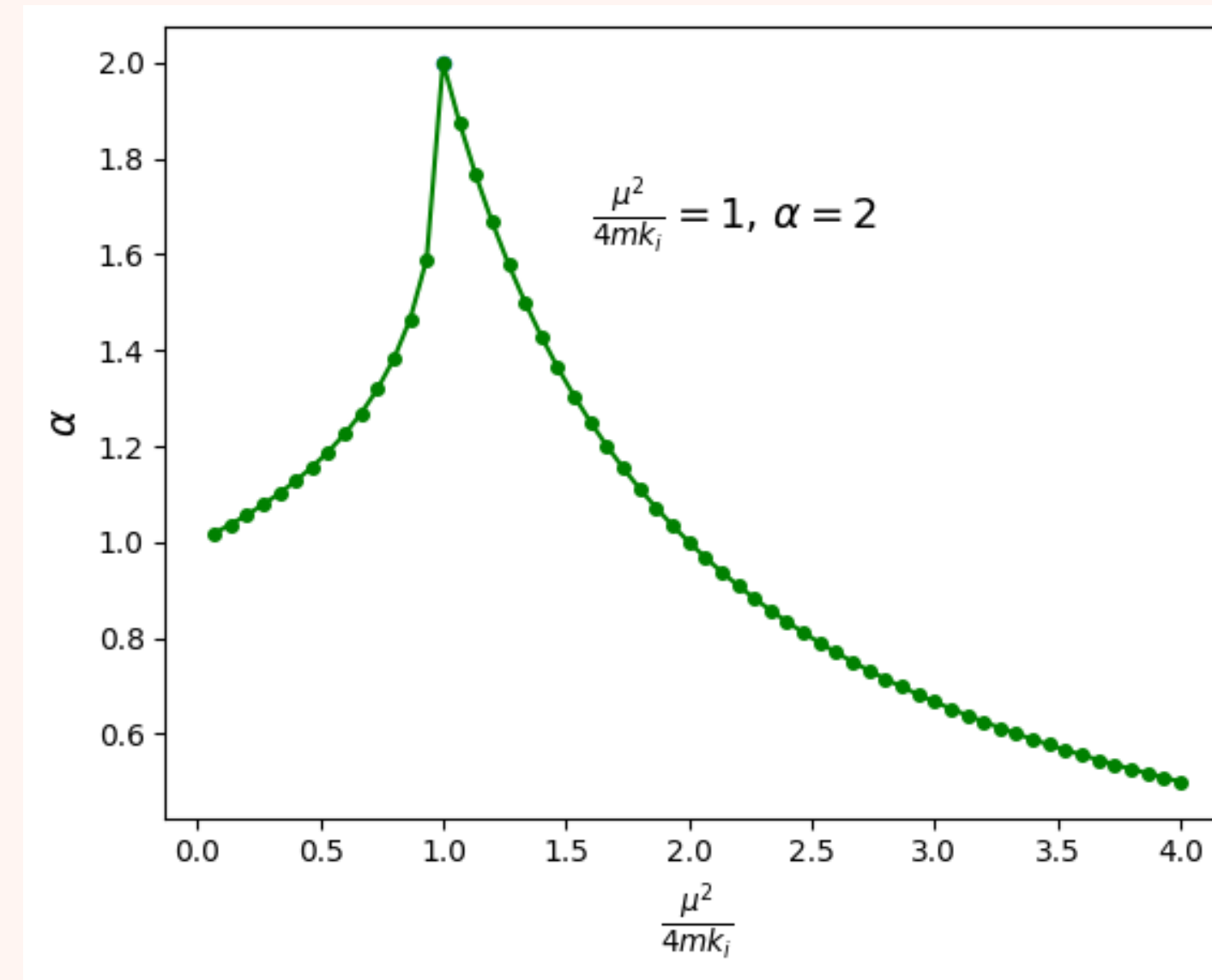


Figure 2. Speed-up provided by momentum

so momentum is most effective when $k_i = \mu^2/4m$, where speed up is doubled

Discrete Case

On discretization and approximating the gradient using the Hessian matrix,

$$w'_{i,t+1} = [(1+p)\mathcal{I} - \epsilon k_i] w'_{i,t} - p w'_{i,t-1} \quad (10)$$

Adding a dummy equation $w'_{i,t} = w'_{i,t}$, we get the following

$$\begin{pmatrix} w'_{i,t} \\ w'_{i,t+1} \end{pmatrix} = A^t \begin{pmatrix} w'_{i,0} \\ w'_{i,1} \end{pmatrix} \quad (11)$$

where $A = \begin{pmatrix} 0 & 1 \\ -p & 1+p-\epsilon k_i \end{pmatrix}$, whose eigen values are given by

$$\lambda_{i,1,2} = \frac{1+p-\epsilon k_i \pm \sqrt{(1+p-\epsilon k_i)^2 - 4p}}{2} \quad (12)$$

Observations:

- For convergence, $\max(|\lambda_{i,1}|, |\lambda_{i,2}|) < 1$ which happens if and only if $-1 < p < 1$ and $0 < \epsilon k_i < 2 + 2p$
- Without $p = 0$, we have $0 < \epsilon k_i < 2$, thus momentum extends range of allowed epsilon, nearly doubles it.
- Compared to the continuous case, discrete system is guaranteed to converge only for some ϵ, k_i and p .
- For small ϵk_i , optimal momentum $p = (1 - \sqrt{\epsilon k_i})^2$ and corresponding $\lambda_{i,1,2} = 1 - \sqrt{\epsilon k_i} < 1 - \epsilon k_i = \lambda_{i,0}$

Experiments

Consider an error function $\mathbf{w}^T A \mathbf{w}$ with $A = \begin{pmatrix} 3 & 0 \\ 0 & 30 \end{pmatrix}$. In the continuous time limit with $\mu = 4, m = 2$, and using large time step $\Delta t = 0.5$ in the discrete case

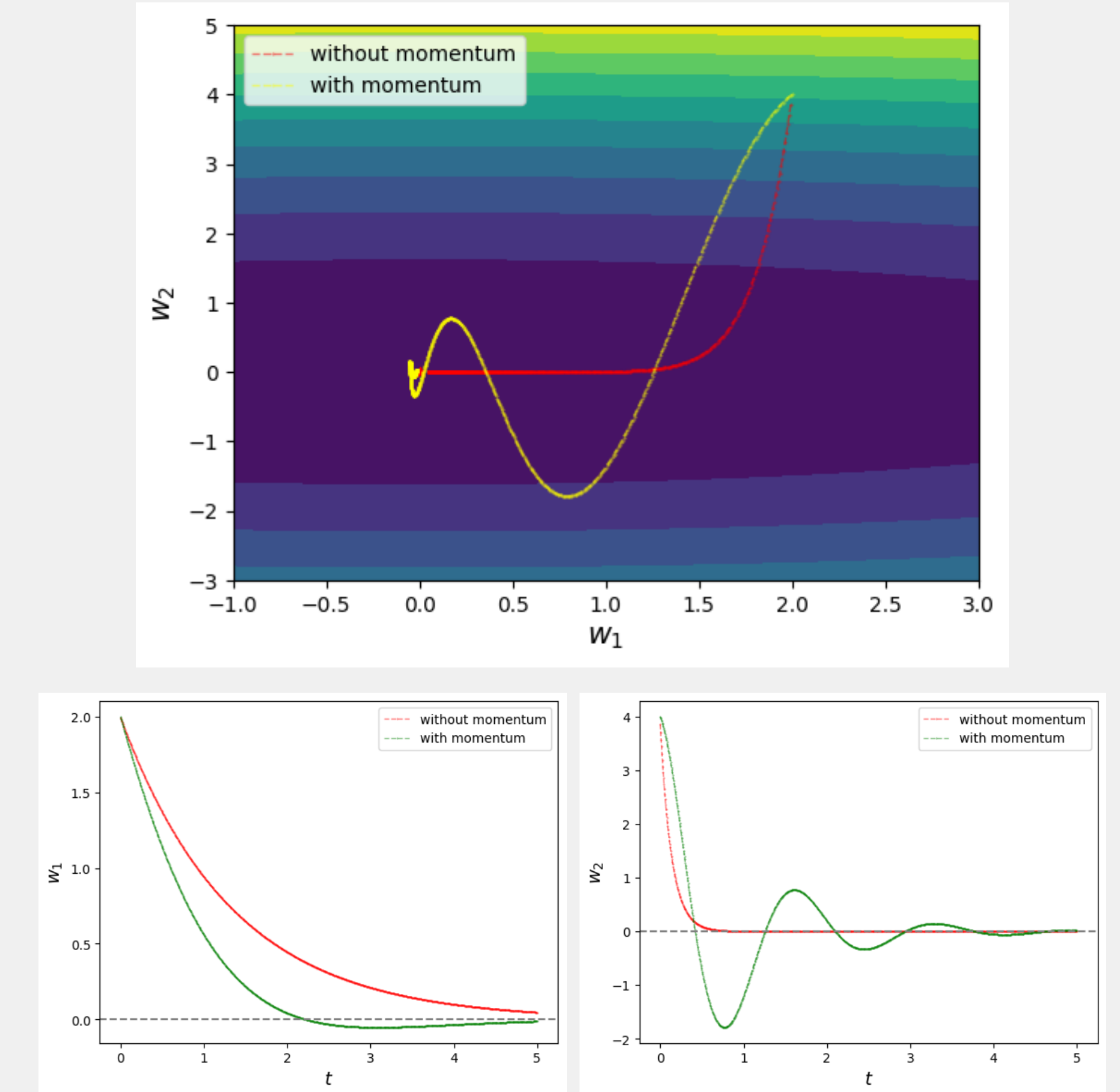


Figure 3. Continuous time

momentum improves convergence along the some directions.

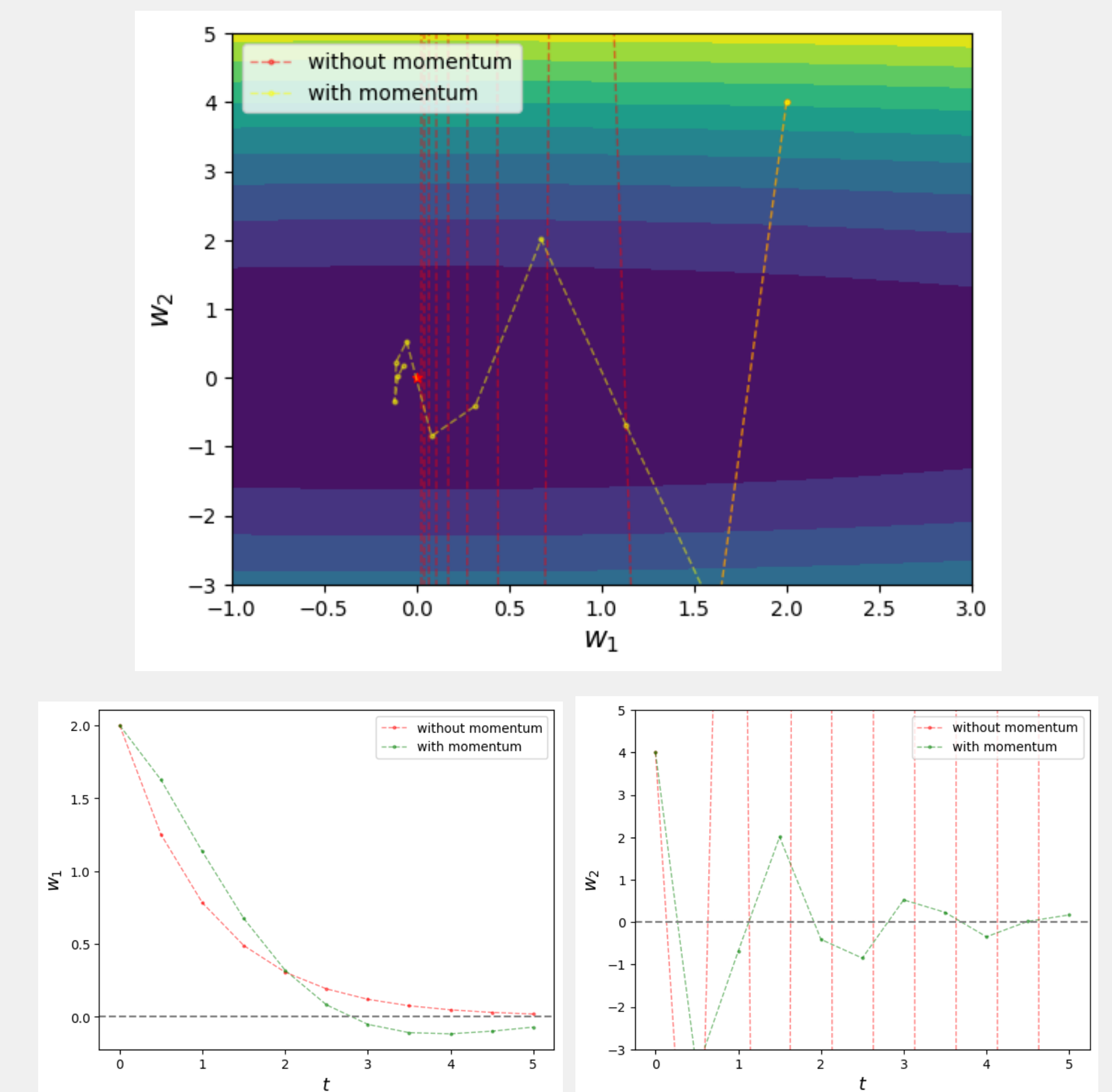


Figure 4. Discrete time

momentum allows convergence with large time steps as expected.

References

- [1] Ning Qian.
On the momentum term in gradient descent learning algorithms.
Neural Networks, 12(1):145-151, 1999.