

# 2-4 Gradient Descent Algorithms and Their Variations

Zhonglei Wang

WISE and SOE, XMU, 2025

# Contents

1. Gradient descent algorithms

2. Variation

3. Learning rate decay

# Recall

## 1. Batch gradient descent

- Randomly initialize  $\boldsymbol{\theta}^{(0)}$
- Based on the current parameter  $\boldsymbol{\theta}^{(t)}$ , obtain

$$\nabla \mathcal{J}(\boldsymbol{\theta}^{(t)}) = \frac{\partial \mathcal{J}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}^{(t)}) = n^{-1} \sum_{i=1}^n \frac{\partial \mathcal{L}_i}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}^{(t)})$$

▷  $\mathcal{L}_i$  : loss function associated with the  $i$ th training example

- Update parameter

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \alpha \nabla \mathcal{J}(\boldsymbol{\theta}^{(t)})$$

- Go back to Step 2 until convergence

# Recall

## 1. Disadvantages

- Computationally inefficient when  $n$  is large

## 2. Why not use part of the training examples for each iteration?

# Mini-batch gradient descent

1. Idea: use part of training examples for each iteration
2. Partition the index set of training examples:  $\{1, \dots, n\} = S_1 \cup \dots \cup S_k$ 
  - $|S_1| = \dots = |S_{k-1}| = m$
  - $|S_k| \leq m$
  - $k = \lceil n/m \rceil$
  - $m$  is usually a power of 2, in practice. For example,  $m = 512$

# Mini-batch gradient descent

1. For the  $t$ th iteration, update the model parameter by

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \alpha \nabla \mathcal{J}_t \left( \boldsymbol{\theta}^{(t)} \right)$$

$$\nabla \mathcal{J}_t \left( \boldsymbol{\theta}^{(t)} \right) = |S_{t \% k + 1}|^{-1} \sum_{i \in S_{t \% k + 1}} \frac{\partial \mathcal{L}_i}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}^{(t)})$$

- Only use training examples in  $S_{t \% k + 1}$  to obtain the gradients

2. We finish one **epoch** when each training example is used

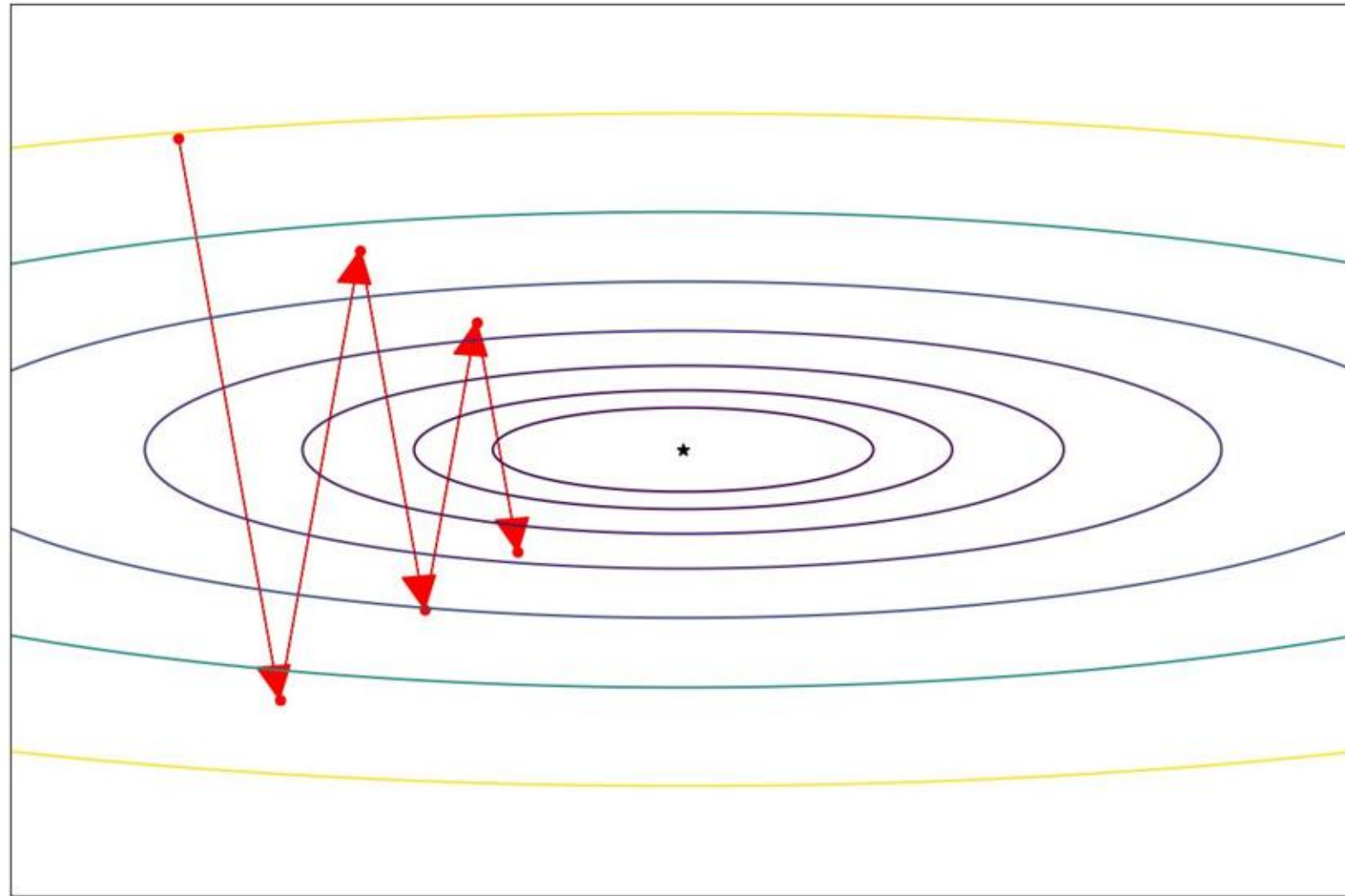
# Discussion

1. Two special cases
  - $m = 1$  : stochastic gradient descent
  - $m = n$  : batch gradient descent
2. Mini-batch gradient descent sacrifices accuracy when  $m < n$
3. Nevertheless, it saves memory and is computationally more efficient
4. Commonly used in practice
5. Next, we consider some efficient gradient descent algorithms.



# Background

1. Gradient shows the fastest direction along which the cost function increases
2. It may causes problems, especially when the cost function is unstable

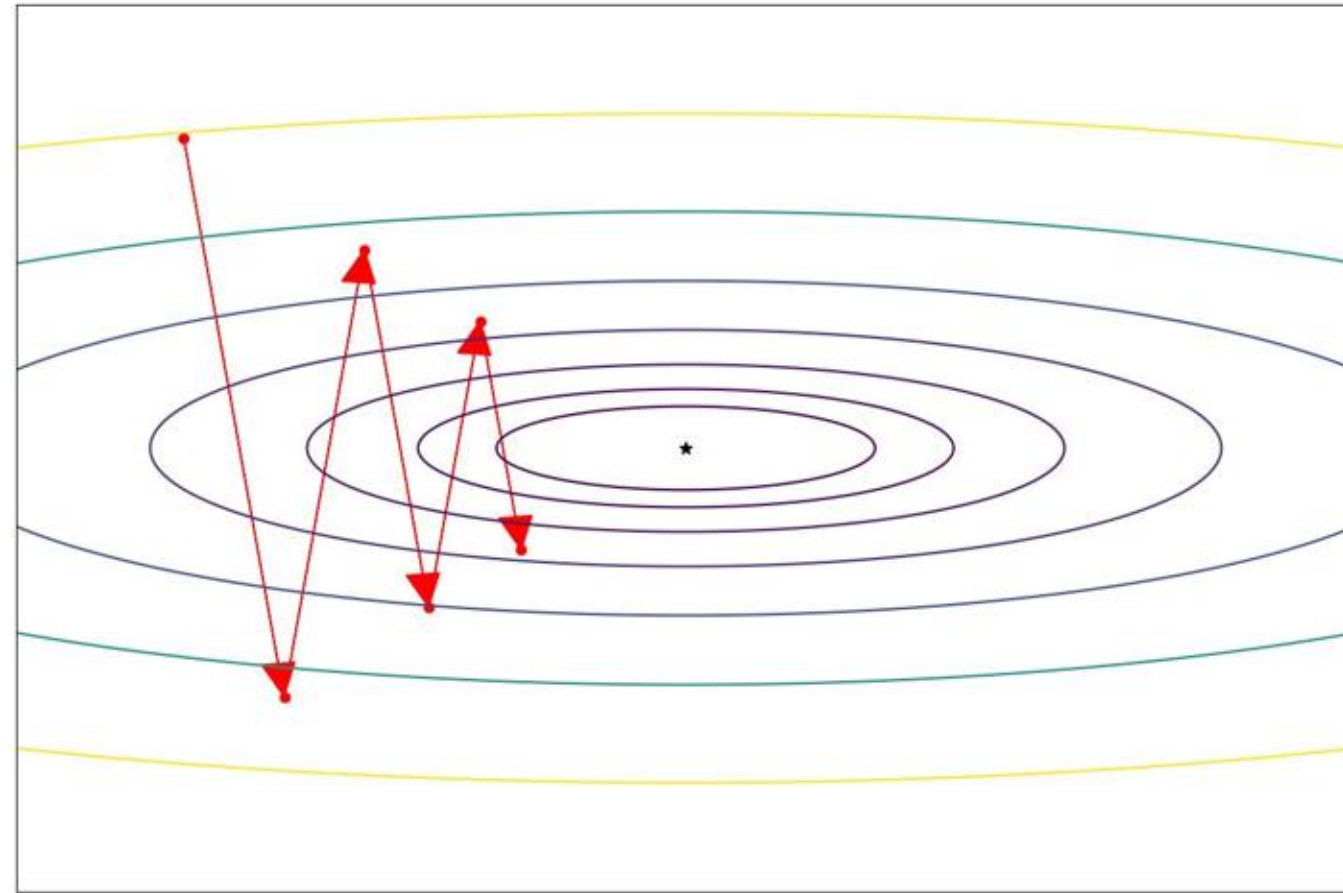




# Background

1. For gradient descent algorithms, replace gradients by possibly more efficient vectors
2. To achieve this goal, we introduce
  - Momentum
  - RMSprop
  - Adam

# Momentum



1. Ideally, we want to

- Decrease the up-down effect
- Increase the left-right effect

# Momentum

## 1. Taking mean might be good?

- Requires memory to store all the past gradients
- Impossible for deep learning models

## 2. Consider EWMA (Exponential Weighted Moving Average)

- Original sequence:  $\{s_i : i = 1, 2, \dots\}$
- EWMA sequence:  $\{v_i : i = 1, 2, \dots\}$

$$v_i = \beta_1 v_{i-1} + (1 - \beta_1) s_i \quad (i = 1, 2, \dots)$$

▷  $v_0 = 0$

▷  $\beta_1$  controls how much we stay with the “momentum”  $v_{i-1}$

# Computation

1. Let  $\mathbf{g}^{(t)}$  be a gradient evaluated at the current parameter for the  $t$ th iteration
  - It can be  $d\mathbf{b}$  or  $d\mathbf{W}$  evaluated at the current parameter  $\boldsymbol{\theta}^{(t)}$
  - We ignore the superscript for layer for simplicity

2. Obtain the EWMA “gradient”  $\mathbf{v}_b^{(t)}$  as follows

$$\mathbf{v}_b^{(t)} = \beta_1 \mathbf{v}_b^{(t-1)} + (1 - \beta_1) \mathbf{g}^{(t)}$$

- $\beta_1 = 0.9$ : a hyperparameter, but seldom tuned
- $\mathbf{v}_b^{(0)} = 0$ : initial momentum

# Computation

## 1. A fact

$$\frac{1}{(1 - \beta_1)^{-1}} \sum_{i=1}^{\infty} \beta_1^i = 1$$

## 2. More details

$$\begin{aligned} \mathbf{v}^{(t)} &= \beta_1 \mathbf{v}^{(t-1)} + (1 - \beta_1) \mathbf{g}^{(t)} \\ &= (1 - \beta_1) \beta_1^{t-1} \mathbf{g}^{(1)} + (1 - \beta_1) \beta_1^{t-2} \mathbf{g}^{(2)} + \dots + (1 - \beta_1) \mathbf{g}^{(t)} \\ &= \frac{1}{(1 - \beta_1)^{-1}} \sum_{i=1}^t \beta_1^{t-i} \mathbf{g}^{(i)} \end{aligned}$$

- When  $t$  is large, it is approximately weighted average
- $(1 - \beta_1)^{-1}$  : can be viewed as the “effective sample size” for EMWA



# Momentum-based gradient descent algorithm

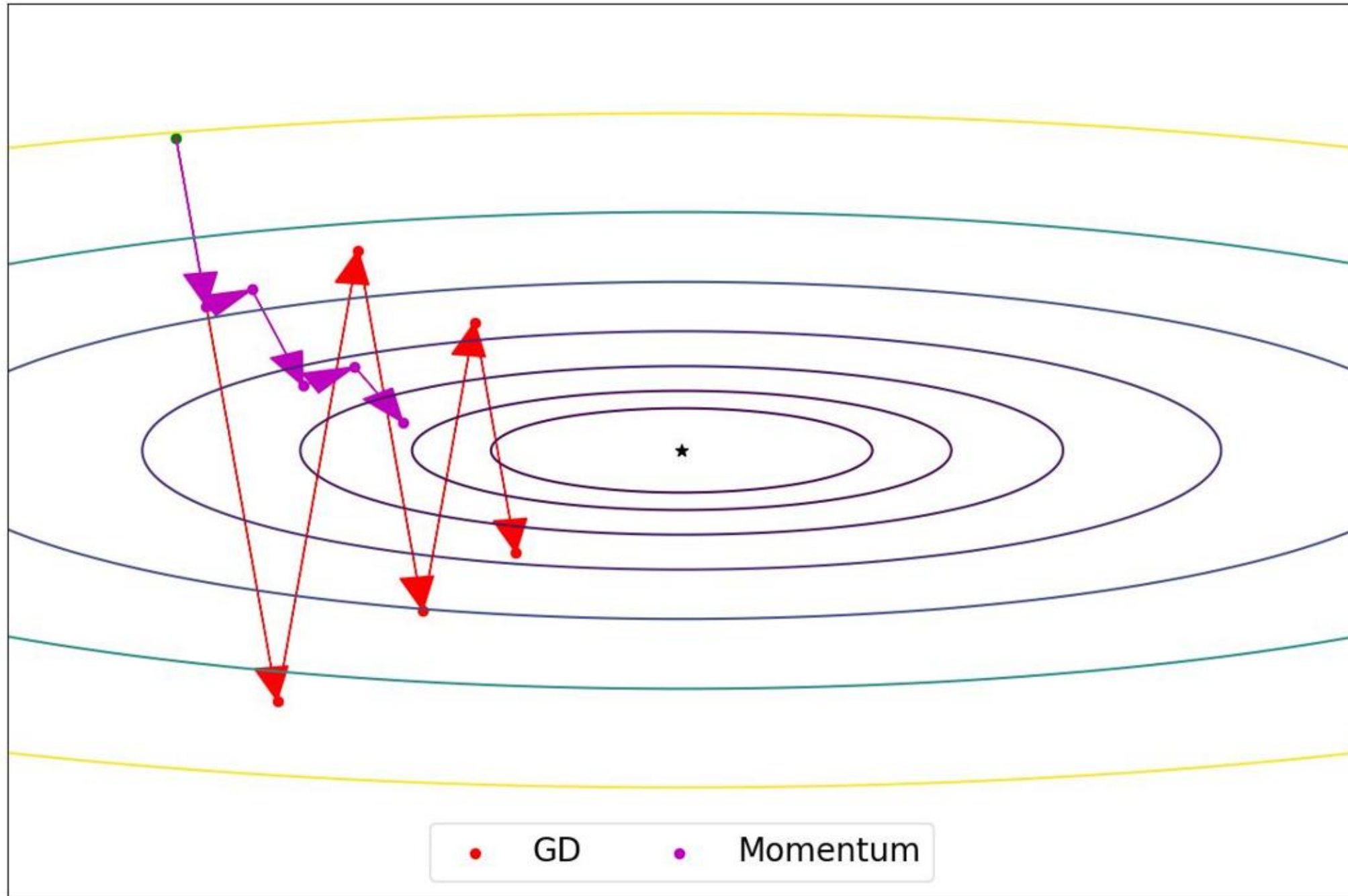
1. Randomly initialize  $\theta^{(0)}$
2. Based on the current model parameter  $\theta^{(t)}$ , obtain  $d\mathbf{b}^{[1](t)}$ 
  - Take the update procedure for  $\mathbf{b}^{[1]}$  as an example
  - The procedure applies to other parameters as well

## 3. Update

$$\mathbf{b}^{[1](t+1)} = \mathbf{b}^{[1](t)} - \alpha \mathbf{v}_b^{[1](t+1)}$$

- $\mathbf{v}_b^{[1](t+1)} = \beta_1 \mathbf{v}_b^{[1](t)} + (1 - \beta_1) d\mathbf{b}^{[1](t)}$
  - $\mathbf{v}_b^{[1](0)} = 0$
  - $\beta_1 = 0.9$  by default
4. Go back to Step 2 until convergence

# Comparison





# Comparison

1. Momentum indeed decreases the variability in the up-down direction
2. Nevertheless, the convergence rate is slow
3. One possible solution: different learning rates for different directions
  - **Low** learning rate for the **up-down** direction
  - **High** learning rate for the **left-right** direction

# RMSPprop-based gradient descent algorithm

Step 1. Randomly initialize  $\theta^{(0)}$

Step 2. Based on the current model parameter  $\theta^{(t)}$ , obtain  $d\mathbf{b}^{[1](t)}$

- Take the update procedure for  $\mathbf{b}^{[1]}$  as an example

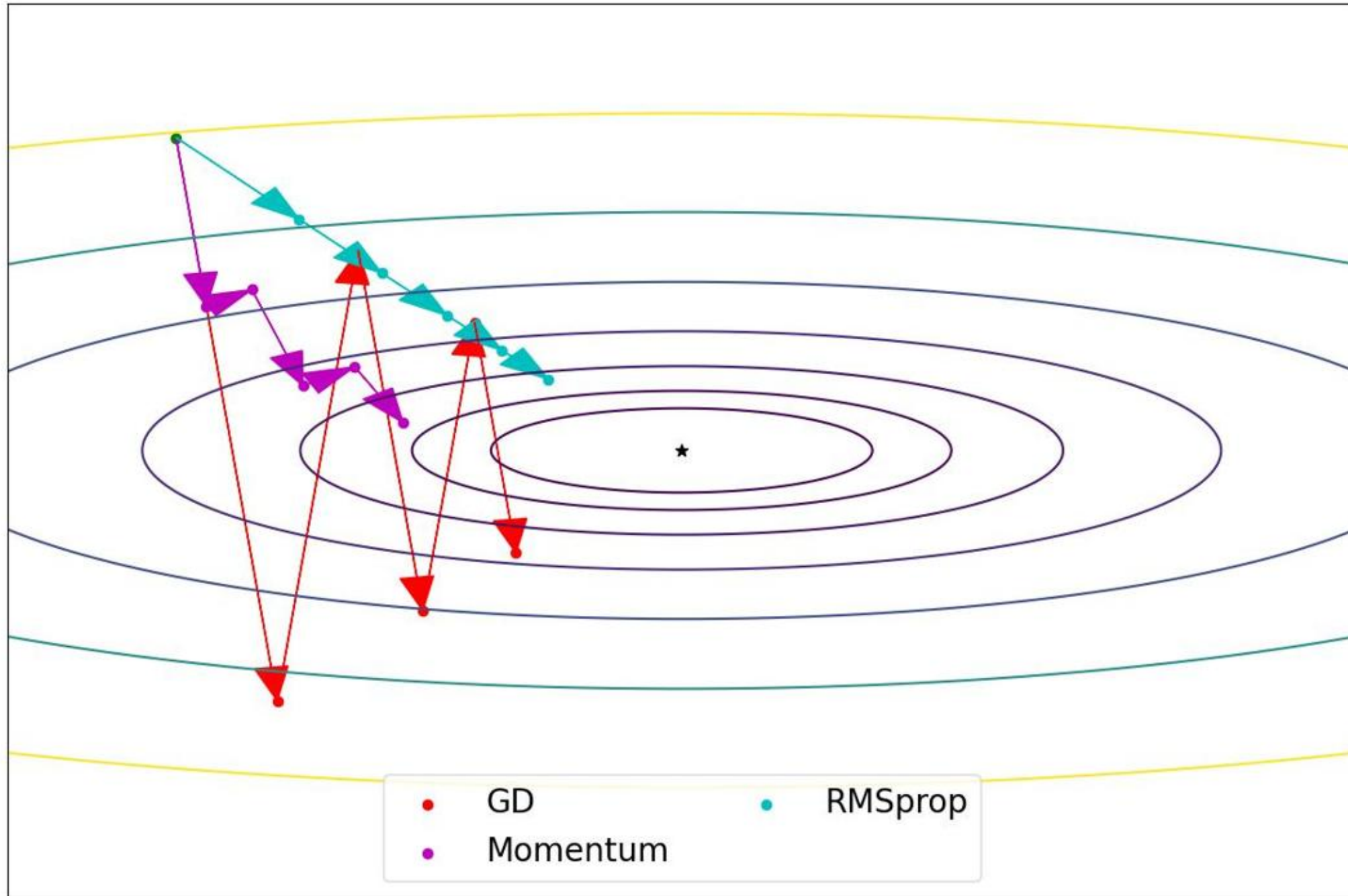
Step 3. Update

$$\mathbf{b}^{[1](t+1)} = \mathbf{b}^{[1](t)} - \frac{\alpha}{\sqrt{\epsilon + \mathbf{s}_b^{[1](t+1)}}} d\mathbf{b}^{[1](t)}$$

- $\epsilon = 10^{-8}$  by default
- $\mathbf{s}_b^{[1](t+1)} = \beta_2 \mathbf{s}_b^{[1](t)} + (1 - \beta_2) d\mathbf{b}^{[1](t)} \circ d\mathbf{b}^{[1](t)}$
- $\mathbf{s}_b^{[1](0)} = 0$
- $\beta_2 = 0.99$  by default

Step 4. Go back to Step 2 until convergence

# Comparison



# Comparison

1. The path for RMSprop is more smooth
2. Nevertheless, the convergence rate is not that fast

# Algorithms

1. Classical gradient descent algorithms

$$\mathbf{b}^{[1]}(t+1) = \mathbf{b}^{[1]}(t) - \alpha d\mathbf{b}^{[1]}(t)$$

2. Momentum only modifies the gradient part

$$\mathbf{b}^{[1]}(t+1) = \mathbf{b}^{[1]}(t) - \alpha \mathbf{v}_b^{[1]}(t+1)$$

3. RMSprop only modifies the learning rate part

$$\mathbf{b}^{[1]}(t+1) = \mathbf{b}^{[1]}(t) - \frac{\alpha}{\sqrt{\epsilon + \mathbf{s}_b^{[1]}(t+1)}} d\mathbf{b}^{[1]}(t)$$

4. Why not combine those two together?



# Adam (Adaptive moment estimation)

1. Randomly initialize  $\theta^{(0)}$
2. Based on the current model parameter  $\theta^{(t)}$ , obtain  $d\mathbf{b}^{[1]}(t)$ 
  - Take the update procedure for  $\mathbf{b}^{[1]}$  as an example

## 3. Update

$$\mathbf{b}^{[1]}(t+1) = \mathbf{b}^{[1]}(t) - \frac{\alpha}{\epsilon + \sqrt{\hat{\mathbf{s}}_b^{[1]}(t+1)}} \hat{\mathbf{v}}_b^{[1]}(t+1)$$

- See next slide for  $\hat{\mathbf{s}}_b^{[1]}(t+1)$  and  $\hat{\mathbf{v}}_b^{[1]}(t+1)$
  - $\epsilon = 10^{-8}$  by default
4. Go back to Step 2 until convergence

# Adam (Adaptive moment estimation)

1. Computation details for  $\hat{\mathbf{s}}_b^{[1](t+1)}$  and  $\hat{\mathbf{v}}_b^{[1](t+1)}$

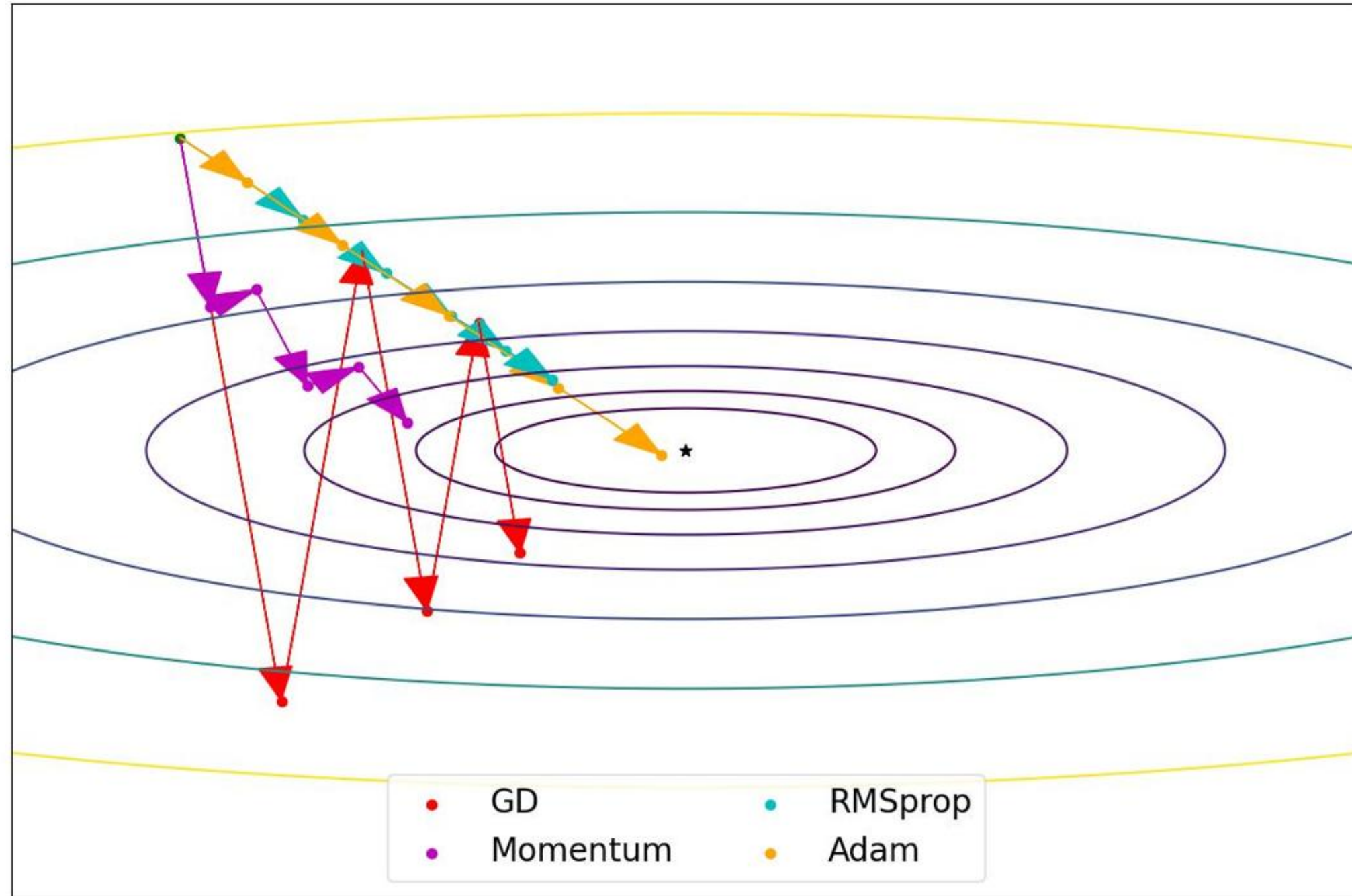
$$\mathbf{s}_b^{[1](t+1)} = \beta_2 \mathbf{s}_b^{[1](t)} + (1 - \beta_2) d\mathbf{b}^{[1](t)} \quad \hat{\mathbf{s}}_b^{[1](t+1)} = \frac{\mathbf{s}_b^{[1](t+1)}}{1 - \beta_2^{t+1}}$$

$$\mathbf{v}_b^{[1](t+1)} = \beta_1 \mathbf{v}_b^{[1](t)} + (1 - \beta_1) d\mathbf{b}^{[1](t)} \quad \hat{\mathbf{v}}_b^{[1](t+1)} = \frac{\mathbf{v}_b^{[1](t+1)}}{1 - \beta_1^{t+1}}$$

- $\mathbf{v}_b^{[1](0)} = \mathbf{s}_b^{[1](0)} = 0$
- $\beta_1 = 0.9$  by default
- $\beta_2 = 0.99$  by default



# Comparison



# Comparison

1. Adam performs the best
2. Adam is commonly used in practice

# Learning rate decay

## 1. Intuition

- As iteration goes, the parameters should get closer to the theoretical values
- It is inefficient to use the SAME learning rate for all iterations
  - ▷ [Inefficient] Small learning rate guarantees good performance, but it takes longer to get conver
  - ▷ [Unstable] A large learning rate leads to instability when iteration index is large
- Actually, we should decrease the learning rate in a reasonable manner

# Learning rate decay

## 1. Denote

- $t$  : iteration index
- $epoch$  : epoch index

## 2. Several possible ways to decay the learning rate

$$\alpha_t = \frac{\alpha_0}{1 + \gamma \cdot epoch}$$

$$\alpha_t = \frac{\alpha_0}{\sqrt{epoch}}$$

$$\alpha_t = 0.95^{epoch} \cdot \alpha_0$$

$$\alpha_t = \frac{1}{\sqrt{t}} \cdot \alpha_0$$

$$\alpha_t = 0.95^t \cdot \alpha_0$$

- $\alpha_0 = 5 \times 10^{-3}$  for example
- $\gamma = 1$  (by default)