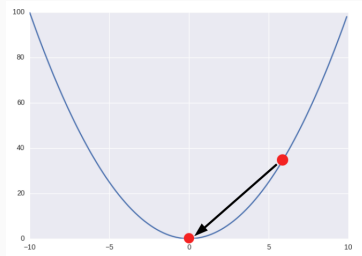


Basic algorithms

Stochastic gradient descent (SGD)

- Most used algorithm for deep learning
- Do not confuse with (deterministic) gradient descent
 - Stochastic uses minibatches
- Algorithm is similar, but there are some important modifications



Gradient descent algorithm

- Full training samples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with targets $\mathbf{y}^{(i)}$
- Compute gradient

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \left(\sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \right) \quad (5)$$

- Apply update

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \mathbf{g} \quad (6)$$

where

- ϵ is the learning rate
- $\boldsymbol{\theta}$ are the network parameters
- $L(\cdot)$ is the loss function

Stochastic gradient descent algorithm

- Minibatch of training samples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with targets $\mathbf{y}^{(i)}$
- Compute gradient

$$\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \left(\sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \right) \quad (7)$$

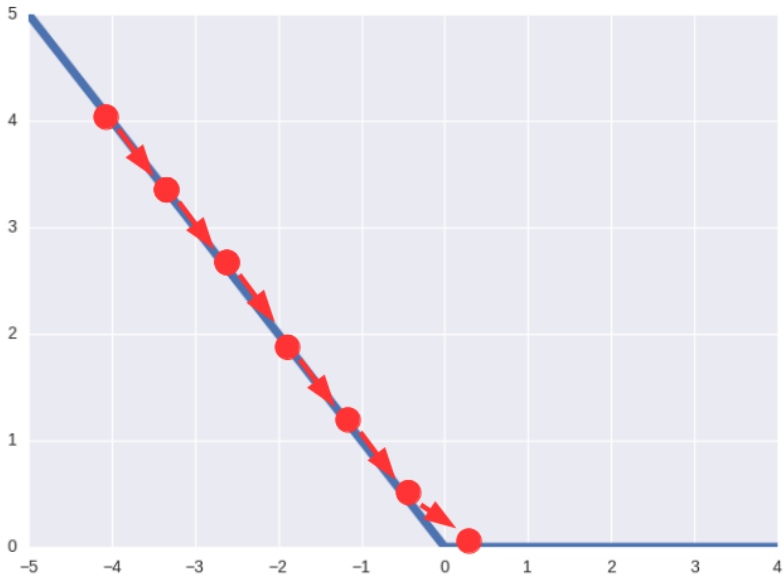
- Apply update

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon_R \hat{\mathbf{g}} \quad (8)$$

Learning rate for SGD

- The main problem is how to choose ϵ_0
- Typically:
 - Higher than the best value for the first 100 iterations
 - Monitor the initial results and use a higher value
 - Too high will cause instability

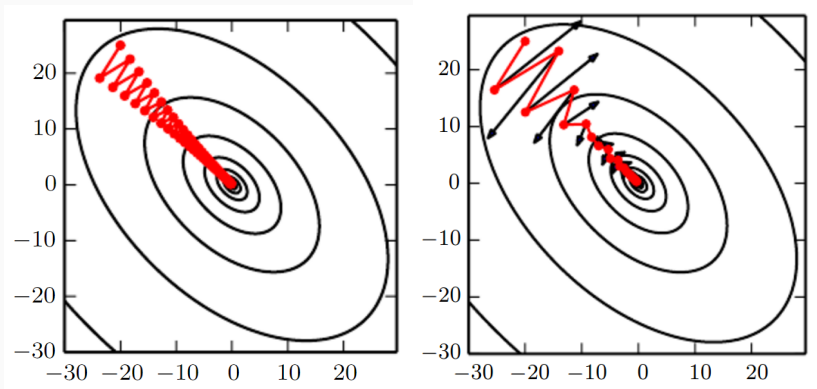
Momentum



Momentum

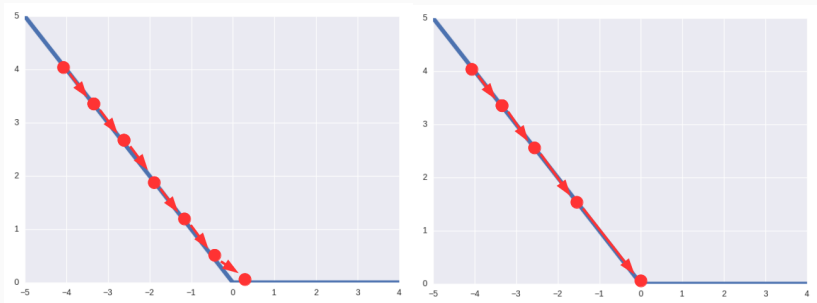
- In these cases, momentum can help
- Derived from the physics term ($= \text{mass} \times \text{velocity}$)
- Assume unit mass, so just consider velocity

Momentum



Source: [Goodfellow et al., 2016]

Momentum



SGD with momentum

- Compute gradient

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \left(\sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \right) \quad (13)$$

- Compute velocity update

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g} \quad (14)$$

- Apply update

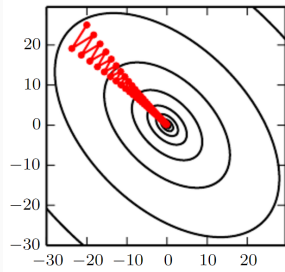
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v} \quad (15)$$

where

- α is the momentum coefficient

Adaptive learning rate

- Learning rate is one of the hyperparameter that impacts the most
- The gradient is highly sensitive to some directions
- If we assume that the sensitivity is axis-aligned, it makes sense to use separate rates for each parameter



Source:
[Goodfellow et al., 2016]

Delta-bar-delta [Jacobs, 1988]

- Early heuristic approach
- Simple idea: if the partial derivative in respect to one parameter remains the same, increase the learning rate, otherwise, decrease
- Must be used in batch methods

AdaGrad [Duchi et al., 2011]

- Scale the gradient according to the historical norms
- Learning rates of parameters with high partial derivatives decrease fast
- Enforces progress in more gently sloped directions
- Nice properties for convex optimization
- But for deep learning decrease the rate in excess

AdaGrad algorithm

- Accumulate squared gradients

$$\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g} \quad (18)$$

- Element-wise update

$$\Delta \boldsymbol{\theta} \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g} \quad (19)$$

- Update parameters

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta} \quad (20)$$

where

- \mathbf{g} is the gradient
- δ is a small constant for stabilization

RMSProp [Hinton, 2012]

- Modification of AdaGrad to perform better on non-convex problems
- AdaGrad accumulates since beginning, gradient may be too small before reaching a convex structure
- RMSProp uses an exponentially weighted moving average

RMSProp algorithm

- Accumulate squared gradients

$$\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g} \quad (21)$$

- Element-wise update

$$\Delta \boldsymbol{\theta} \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g} \quad (22)$$

- Update parameters

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta} \quad (23)$$

where

- ρ is the decay rate

Adam [Kingma and Ba, 2014]

- Adaptive Moments, variation of RMSProp + Momentum
- Momentum is incorporated directly as an estimate of the first order moment
 - In RMSProp momentum is included after rescaling the gradients
- Adam also add bias correction to the moments

Adam algorithm

- Update time step: $t \leftarrow t + 1$
- Update biased moment estimates

$$\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g} \quad (24)$$

$$\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g} \quad (25)$$

- Correct biases

$$\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t} \quad (26)$$

$$\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t} \quad (27)$$

- Update parameters

$$\Delta \boldsymbol{\theta} \leftarrow -\epsilon \frac{\hat{\mathbf{s}}}{\delta + \sqrt{\hat{\mathbf{r}}}} \quad (28)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta} \quad (29)$$