

different scales, so using Momentum optimization helps a lot. It can also help roll past local optima.



Due to the momentum, the optimizer may overshoot a bit, then come back, overshoot again, and oscillate like this many times before stabilizing at the minimum. This is one of the reasons why it is good to have a bit of friction in the system: it gets rid of these oscillations and thus speeds up convergence.

Implementing Momentum optimization in TensorFlow is a no-brainer: just replace the `GradientDescentOptimizer` with the `MomentumOptimizer`, then lie back and profit!

```
optimizer = tf.train.MomentumOptimizer(learning_rate=learning_rate,  
                                         momentum=0.9)
```

The one drawback of Momentum optimization is that it adds yet another hyperparameter to tune. However, the momentum value of 0.9 usually works well in practice and almost always goes faster than Gradient Descent.

## Nesterov Accelerated Gradient

One small variant to Momentum optimization, proposed by **Yurii Nesterov in 1983**,<sup>12</sup> is almost always faster than vanilla Momentum optimization. The idea of *Nesterov Momentum optimization*, or *Nesterov Accelerated Gradient* (NAG), is to measure the gradient of the cost function not at the local position but slightly ahead in the direction of the momentum (see **Equation 11-5**). The only difference from vanilla Momentum optimization is that the gradient is measured at  $\theta + \beta \mathbf{m}$  rather than at  $\theta$ .

*Equation 11-5. Nesterov Accelerated Gradient algorithm*

1.  $\mathbf{m} \leftarrow \beta \mathbf{m} + \eta \nabla_{\theta} J(\theta + \beta \mathbf{m})$
2.  $\theta \leftarrow \theta - \mathbf{m}$

This small tweak works because in general the momentum vector will be pointing in the right direction (i.e., toward the optimum), so it will be slightly more accurate to use the gradient measured a bit farther in that direction rather than using the gradient at the original position, as you can see in **Figure 11-6** (where  $\nabla_1$  represents the gradient of the cost function measured at the starting point  $\theta$ , and  $\nabla_2$  represents the gradient at the point located at  $\theta + \beta \mathbf{m}$ ). As you can see, the Nesterov update ends up

<sup>12</sup> “A Method for Unconstrained Convex Minimization Problem with the Rate of Convergence  $O(1/k^2)$ ,” Yurii Nesterov (1983).

slightly closer to the optimum. After a while, these small improvements add up and NAG ends up being significantly faster than regular Momentum optimization. Moreover, note that when the momentum pushes the weights across a valley,  $\nabla_1$  continues to push further across the valley, while  $\nabla_2$  pushes back toward the bottom of the valley. This helps reduce oscillations and thus converges faster.

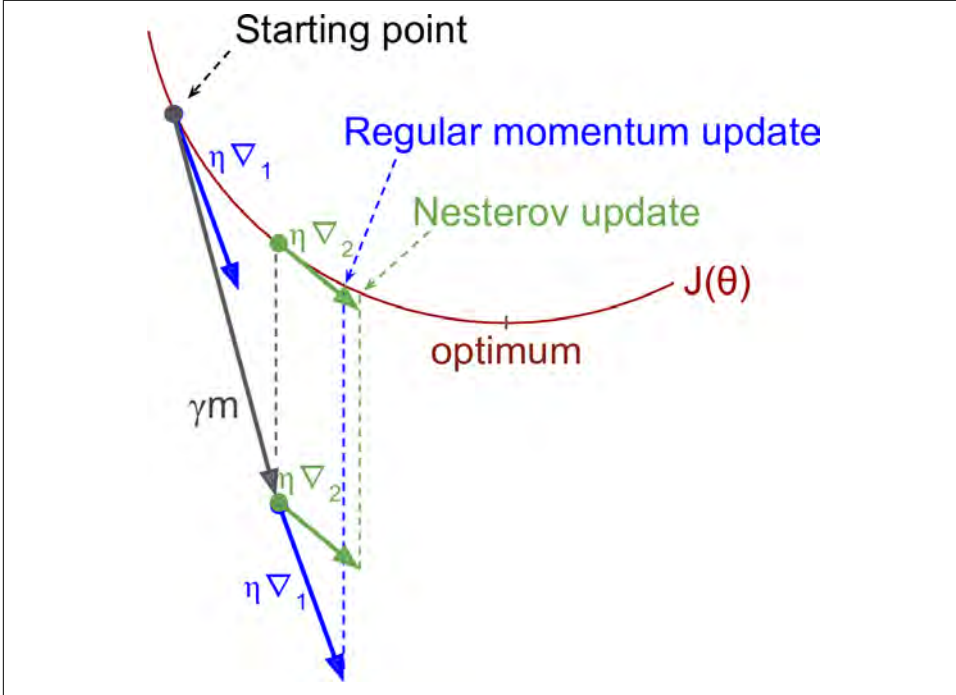


Figure 11-6. Regular versus Nesterov Momentum optimization

NAG will almost always speed up training compared to regular Momentum optimization. To use it, simply set `use_nesterov=True` when creating the `MomentumOptimizer`:

```
optimizer = tf.train.MomentumOptimizer(learning_rate=learning_rate,
                                         momentum=0.9, use_nesterov=True)
```

## AdaGrad

Consider the elongated bowl problem again: Gradient Descent starts by quickly going down the steepest slope, then slowly goes down the bottom of the valley. It would be nice if the algorithm could detect this early on and correct its direction to point a bit more toward the global optimum.