

---

**Algorithm A.10 Simple gradient ascent algorithm**


---

```

Procedure Gradient-Ascent (
     $\theta^1$ , // Initial starting point
     $f_{\text{obj}}$ , // Function to be optimized
     $\delta$  // Convergence threshold
)
1    $t \leftarrow 1$ 
2   do
3      $\theta^{t+1} \leftarrow \theta^t + \eta \nabla f_{\text{obj}}(\theta^t)$ 
4      $t \leftarrow t + 1$ 
5   while  $\|\theta^t - \theta^{t-1}\| > \delta$ 
6   return  $(\theta^t)$ 

```

---

search of algorithm A.5 (see appendix A.4.2). Using the Taylor expansion of a function, we know that, in the neighborhood of  $\theta^0$ , the function can be approximated by the linear equation

$$f_{\text{obj}}(\theta) \approx f_{\text{obj}}(\theta^0) + (\theta - \theta^0)^T \nabla f_{\text{obj}}(\theta^0).$$

Using basic properties of linear algebra, we can check that the slope of this linear function, that is,  $\nabla f_{\text{obj}}(\theta^0)$ , points to the direction of the steepest ascent. This observation suggests that, if we take a step in the direction of the gradient, we increase the value of  $f_{\text{obj}}$ . This reasoning leads to the simple gradient ascent algorithm shown in algorithm A.10. Here,  $\eta$  is a constant that determines the *rate* of ascent at each iteration. Since the gradient  $\nabla f_{\text{obj}}$  approaches 0 as we approach a maximum point, the procedure will converge if  $\eta$  is sufficiently small.

Note that, in order to apply gradient ascent, we need to be able to evaluate the function  $f_{\text{obj}}$  at different points, and also to evaluate its gradient. In several examples we encounter in this book, we can perform these calculations, although in some cases these are costly. Thus, a major objective is to reduce the number of points at which we evaluate  $f_{\text{obj}}$  or  $\nabla f_{\text{obj}}$ .

The performance of gradient ascent depends on the choice of  $\eta$ . If  $\eta$  is too large, then the algorithm can “overshoot” the maximum in each iteration. For sufficiently small value of  $\eta$ , the gradient ascent algorithm will converge, but if  $\eta$  is too small, we will need many iterations to converge. Thus, one of the difficult points in applying this algorithm is deciding on the value of  $\eta$ . Indeed, in practice, one typically needs to begin with a large  $\eta$ , and decrease it over time; this approach leaves us with the problem of choosing an appropriate schedule for shrinking  $\eta$ .

### A.5.2.2 Line Search

An alternative approach is to adaptively choose the step size  $\eta$  at each step. The intuition is that we choose a direction to climb and continue in that direction until we reach a point where we start to descend. In this procedure, at each point  $\theta^t$  in the search, we define a “line” in the direction of the gradient:

$$g(\eta) = \vec{\theta}^t + \eta \nabla f_{\text{obj}}(\theta^t).$$

line search

We now use a *line search* procedure to find the value of  $\eta$  that defines a (local) maximum of