



Figure 4.3: Optimization algorithms may fail to find a global minimum when there are multiple local minima or plateaus present. In the context of deep learning, we generally accept such solutions even though they are not truly minimal, so long as they correspond to significantly low values of the cost function.

critical points are points where every element of the gradient is equal to zero.

The **directional derivative** in direction  $\mathbf{u}$  (a unit vector) is the slope of the function  $f$  in direction  $\mathbf{u}$ . In other words, the directional derivative is the derivative of the function  $f(\mathbf{x} + \alpha\mathbf{u})$  with respect to  $\alpha$ , evaluated at  $\alpha = 0$ . Using the chain rule, we can see that  $\frac{\partial}{\partial \alpha} f(\mathbf{x} + \alpha\mathbf{u})$  evaluates to  $\mathbf{u}^\top \nabla_{\mathbf{x}} f(\mathbf{x})$  when  $\alpha = 0$ .

To minimize  $f$ , we would like to find the direction in which  $f$  decreases the fastest. We can do this using the directional derivative:

$$\min_{\mathbf{u}, \mathbf{u}^\top \mathbf{u} = 1} \mathbf{u}^\top \nabla_{\mathbf{x}} f(\mathbf{x}) \quad (4.3)$$

$$= \min_{\mathbf{u}, \mathbf{u}^\top \mathbf{u} = 1} \|\mathbf{u}\|_2 \|\nabla_{\mathbf{x}} f(\mathbf{x})\|_2 \cos \theta \quad (4.4)$$

where  $\theta$  is the angle between  $\mathbf{u}$  and the gradient. Substituting in  $\|\mathbf{u}\|_2 = 1$  and ignoring factors that do not depend on  $\mathbf{u}$ , this simplifies to  $\min_{\mathbf{u}} \cos \theta$ . This is minimized when  $\mathbf{u}$  points in the opposite direction as the gradient. In other words, the gradient points directly uphill, and the negative gradient points directly downhill. We can decrease  $f$  by moving in the direction of the negative gradient. This is known as the **method of steepest descent** or **gradient descent**.

Steepest descent proposes a new point

$$\mathbf{x}' = \mathbf{x} - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x}) \quad (4.5)$$