

Nabla Operator: Gradients in \mathbb{R}^n are conveniently expressed with the *Nabla operator*:

$$\nabla \equiv \begin{pmatrix} \partial/\partial x_1 \\ \partial/\partial x_2 \\ \vdots \\ \partial/\partial x_n \end{pmatrix} \quad \text{or applied to some } f(\mathbf{x}) \quad \nabla f(\mathbf{x}) = \begin{pmatrix} \partial f/\partial x_1 \\ \partial f/\partial x_2 \\ \vdots \\ \partial f/\partial x_n \end{pmatrix}$$

Hessian Matrix: The Hessian is a square matrix of 2nd order partial derivatives of a scalar-valued function and a measure of local curvature:

$$H_f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

Gradient Descent

Gradient Descent (GD) is a first-order iterative optimization algorithm for finding a local minimum of a differentiable function:

- Assume $F(\mathbf{x})$ is defined and locally differentiable around \mathbf{a}
- Then $F(\mathbf{x})$ decreases fastest in the direction of the *negative* gradient of F at \mathbf{a} , which is $-\nabla F(\mathbf{a})$
- For a sufficiently small step size (learning rate) $\gamma \in \mathbb{R}^+$, we have

$$\mathbf{a}_{n+1} = \mathbf{a}_n - \gamma \nabla F(\mathbf{a}_n) \quad \text{and} \quad F(\mathbf{a}_n) \geq F(\mathbf{a}_{n+1})$$

Iterating the above we get the

Gradient Descent Update Rule:

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \gamma_n \nabla F(\mathbf{x}_n), \quad n \geq 0$$

where $F(\mathbf{x}_0) \geq F(\mathbf{x}_1) \geq F(\mathbf{x}_2) \geq \dots$ is a monotonic sequence

