

Deep Learning

3.3 Gradient Descent

Dr. Konda Reddy Mopuri
kmopuri@iittp.ac.in
Dept. of CSE, IIT Tirupati
Aug-Dec 2021

Training an ML model

- ① Finding the parameters that minimize the training loss

$$W^*, \mathbf{b}^* = \underset{W, \mathbf{b}}{\operatorname{argmin}} \mathcal{L}(f(\cdot; W, \mathbf{b}); \mathcal{D})$$

Training an ML model

- ① Finding the parameters that minimize the training loss

$$W^*, \mathbf{b}^* = \underset{W, \mathbf{b}}{\operatorname{argmin}} \mathcal{L}(f(\cdot; W, \mathbf{b}); \mathcal{D})$$

- ② How do we find these optimal parameters?
 - ① Closed form solution (e.g. linear regression)
 - ② Ad-hoc recipes (e.g. Perceptron, K-NN classifier)

Training an ML model

- ① Finding the parameters that minimize the training loss

$$W^*, \mathbf{b}^* = \underset{W, \mathbf{b}}{\operatorname{argmin}} \mathcal{L}(f(\cdot; W, \mathbf{b}); \mathcal{D})$$

- ② How do we find these optimal parameters?
 - ① Closed form solution (e.g. linear regression)
 - ② Ad-hoc recipes (e.g. Perceptron, K-NN classifier)
 - ③ What if the loss function can't be minimized analytically?
- ③ General minimization method used in such cases is the 'Gradient Descent'.

Gradient

- ① Given a function

$$\begin{aligned}f : \mathcal{R}^D &\rightarrow \mathcal{R} \\ x &\rightarrow f(x_1, x_2, \dots, x_D)\end{aligned}$$

Gradient

- ① Given a function

$$f : \mathcal{R}^D \rightarrow \mathcal{R}$$

$$x \rightarrow f(x_1, x_2, \dots, x_D)$$

- ② Its gradient is the mapping

$$\nabla f : \mathcal{R}^D \rightarrow \mathcal{R}^D$$

$$x \rightarrow \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right)$$

Gradient

- ① Given a function

$$\begin{aligned} f : \mathcal{R}^D &\rightarrow \mathcal{R} \\ x &\rightarrow f(x_1, x_2, \dots, x_D) \end{aligned}$$

- ② Its gradient is the mapping

$$\begin{aligned} \nabla f : \mathcal{R}^D &\rightarrow \mathcal{R}^D \\ x &\rightarrow \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right) \end{aligned}$$

- ③ It computes how much each input component influences the value of f locally.

Gradient

- ① Given a function

$$f : \mathcal{R}^D \rightarrow \mathcal{R}$$

$$x \rightarrow f(x_1, x_2, \dots, x_D)$$

- ② Its gradient is the mapping

$$\nabla f : \mathcal{R}^D \rightarrow \mathcal{R}^D$$

$$x \rightarrow \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right)$$

- ③ It computes how much each input component influences the value of f locally.
- ④ The gradient vector is interpreted as the direction and rate of fastest increase.

Gradient Descent in ML

- ① Goal is to minimize the error (or loss): determine the parameters θ that minimize the loss $\mathcal{L}(\theta)$

Gradient Descent in ML

- ① Goal is to minimize the error (or loss): determine the parameters θ that minimize the loss $\mathcal{L}(\theta)$
- ② Gradient points uphill \rightarrow negative of gradient points downhill

Gradient Descent in ML

- ① Goal is to minimize the error (or loss): determine the parameters θ that minimize the loss $\mathcal{L}(\theta)$
- ② Gradient points uphill \rightarrow negative of gradient points downhill
- ③
 - ① Start with an arbitrary initial parameter vector θ_0
 - ② Repeatedly modify it via updating in small steps
 - ③ At each step, modify in the direction that produces steepest descent along the error surface

Gradient Descent in ML

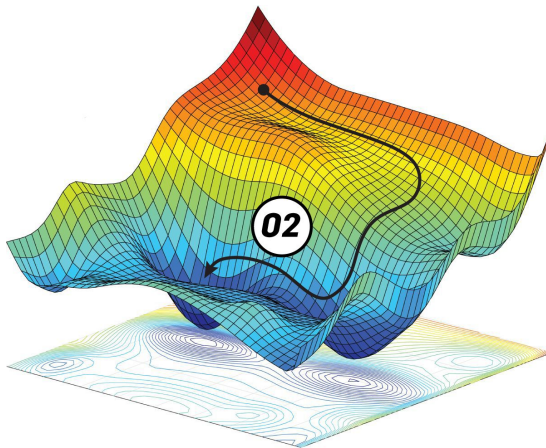


Figure credits:Ahmed Fawzy Gad

Gradient Descent in ML

- ① Start with an arbitrary initial parameter vector θ_0
- ② Repeatedly modify it via updating in small steps
- ③ At each step, modify in the direction that produces steepest descent along the error surface

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t)$$

Gradient Descent in ML

- ① Start with an arbitrary initial parameter vector θ_0
- ② Repeatedly modify it via updating in small steps
- ③ At each step, modify in the direction that produces steepest descent along the error surface

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t)$$

- ④ Almost always ends in a local minimum, choice of parameters θ_0 and η are important.

Gradient descent example

- ① Logistic regression (we will work it out on whiteboard)