

Gradient Descent and Backpropagation in Machine Learning

Jinn-Liang Liu

Institute of Computational and Modeling Science,
National Tsing Hua University, Taiwan
2017/7/11, 2018/12/8

The change of the function (height) $f(\mathbf{x})$ at $\mathbf{x} \in R^2$ in the direction $\mathbf{p} \in R^2$ is a **directional derivative** defined as

$$D_{\mathbf{p}}f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{p}) - f(\mathbf{x})}{t} \quad (1)$$

($\mathbf{x} = (x, y)$ any point, $\mathbf{p} = (p_1, p_2)$ any unit direction)

$$\begin{aligned} &= \lim_{t \rightarrow 0} \frac{f(x + tp_1, y + tp_2) - f(x, y)}{t} \\ &= \lim_{t \rightarrow 0} \left\{ \frac{f(x + tp_1, y + tp_2) - f(x, y + tp_2)}{t} \right. \\ &\quad \left. + \frac{f(x, y + tp_2) - f(x, y)}{t} \right\} \\ &= \lim_{t \rightarrow 0} \left\{ \frac{[f(x + tp_1, y + tp_2) - f(x, y + tp_2)]p_1}{tp_1} \right. \\ &\quad \left. + \frac{[f(x, y + tp_2) - f(x, y)]p_2}{tp_2} \right\} \\ &= \lim_{t \rightarrow 0} \left\{ \frac{[f(x + \Delta x, y + tp_2) - f(x, y + tp_2)]p_1}{\Delta x} \right. \\ &\quad \left. + \frac{[f(x, y + \Delta y) - f(x, y)]p_2}{\Delta y} \right\} \\ &= \frac{\partial f(x, y)}{\partial x} p_1 + \frac{\partial f(x, y)}{\partial y} p_2 \end{aligned}$$

$$(\nabla = \left\langle \frac{\partial}{\partial x}, \frac{\partial}{\partial y} \right\rangle = \text{grad} = \text{del, the gradient operator}) \quad (2)$$

$$= \nabla f(x, y) \cdot \mathbf{p} = |\nabla f| |\mathbf{p}| \cos \theta \quad (3)$$

\Rightarrow Max value of $D_{\mathbf{p}}f$ is $|\nabla f|$ in $\mathbf{p} = \nabla f / |\nabla f|$ with $\theta = 0$

\Rightarrow Min value of $D_{\mathbf{p}}f$ is $-|\nabla f|$ in $\mathbf{p} = -\nabla f / |\nabla f|$ with $\theta = 180^\circ$

Therefore, we should choose the direction $\mathbf{p} = -\nabla f(\mathbf{x})/|\nabla f(\mathbf{x})|$ (negative gradient) if we want to minimize $f(\mathbf{x})$ (to go down to the valley) in the **fastest** way (**steepest descent** or **gradient descent**) when we are at $\mathbf{x} = (x, y)$ on a mountain surface described by the function $z = f(\mathbf{x})$ at the elevation z . The **method of gradient (steepest) descent** is thus an iterative process

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_{k-1}\mathbf{p}_{k-1} \quad (4)$$

of changing (updating) our current location $\mathbf{x}_{k-1} = (x_{k-1}, y_{k-1})$ by deciding the next **stepping length** α_k (called **learning rate** in Machine Learning (ML)) in our **predicted (gradient) direction** $\mathbf{p}_{k-1} = -\nabla f(\mathbf{x}_{k-1})/|\nabla f(\mathbf{x}_{k-1})|$.

Exercise: Show that the **gradient vector** $\nabla f(\mathbf{x}_0)$ is *perpendicular* to the **tangent vector** $\mathbf{r}'(t_0)$ to the level curve $f(\mathbf{x}) = c$ at $\mathbf{x} = \mathbf{x}_0 = (x_0, y_0) = \langle x(t_0), y(t_0) \rangle = \mathbf{r}(t_0)$, where $\mathbf{r}'(t) = \frac{d}{dt} \langle x(t), y(t) \rangle$. Draw a graph with a mountain, surface, level curve, tangent vector, gradient vector, valley, and all math notations.

1D Example: Consider the minimization (optimization) problem

$$\text{Minimize } y = f(x) = x^2, \forall x \in R^1. \quad (5)$$

$$\text{Method: } x_k = x_{k-1} + \alpha_{k-1}p_{k-1} \text{ with } x_0 = 2, \alpha_{k-1} = \frac{1}{2}, \forall k. \quad (6)$$

$$\nabla f(x) = \frac{df(x)}{dx} = 2x \quad (7)$$

$$p_0 = \frac{-\nabla f(x_0)}{|\nabla f(x_0)|} = \frac{-4}{|4|} = -1 \text{ (go west if } x_0 > 0)$$

$$x_1 = x_0 + \alpha_0 p_0 = 2 - \frac{1}{2} = \frac{3}{2} \Rightarrow$$

$$x_2 = x_1 + \alpha_1 p_1 = \frac{3}{2} - \frac{1}{2} = 1 \Rightarrow$$

$$x_3 = x_2 + \alpha_2 p_2 = 1 - \frac{1}{2} = \frac{1}{2} \Rightarrow$$

$$x_4 = x_3 + \alpha_3 p_3 = \frac{1}{2} - \frac{1}{2} = 0 = x^* \text{ (optimizer)} \Rightarrow$$

$$y^* = f(x^*) = 0 \text{ (optimal value).}$$

2D Example: Consider the minimization problem

$$\text{Minimize } z = f(\mathbf{x}) = \frac{x^2}{4^2} + y^2, \forall \mathbf{x} = \langle x, y \rangle \in R^2. \quad (8)$$

$$\text{Method: } \mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_{k-1} \mathbf{p}_{k-1} \text{ with } \mathbf{x}_0 = \left\langle 2, \frac{1}{4} \right\rangle, \alpha_{k-1} = \frac{\sqrt{5}}{4}, \forall k. \quad (9)$$

$$\nabla f(\mathbf{x}) = \left\langle \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right\rangle = \left\langle \frac{x}{8}, 2y \right\rangle \quad (10)$$

$$\mathbf{p}_0 = \frac{-\nabla f(\mathbf{x}_0)}{|\nabla f(\mathbf{x}_0)|} = \frac{-\left\langle \frac{1}{4}, \frac{1}{2} \right\rangle}{\left| \left\langle \frac{1}{4}, \frac{1}{2} \right\rangle \right|} = \frac{-\left\langle \frac{1}{4}, \frac{1}{2} \right\rangle}{\frac{\sqrt{5}}{4}} = -\left\langle \frac{1}{\sqrt{5}}, \frac{2}{\sqrt{5}} \right\rangle$$

$$\mathbf{x}_1 = \mathbf{x}_0 + \alpha_0 \mathbf{p}_0 = \left\langle 2, \frac{1}{4} \right\rangle - \frac{\sqrt{5}}{4} \left\langle \frac{1}{\sqrt{5}}, \frac{2}{\sqrt{5}} \right\rangle = \left\langle \frac{7}{4}, \frac{-1}{4} \right\rangle$$

$$\mathbf{p}_1 = \frac{-\nabla f(\mathbf{x}_1)}{|\nabla f(\mathbf{x}_1)|} = \frac{-\left\langle \frac{7}{32}, \frac{-1}{2} \right\rangle}{\left| \left\langle \frac{7}{32}, \frac{-1}{2} \right\rangle \right|}$$

$$\mathbf{x}_2 = \mathbf{x}_1 + \alpha_1 \mathbf{p}_1 = \dots \Rightarrow \dots$$

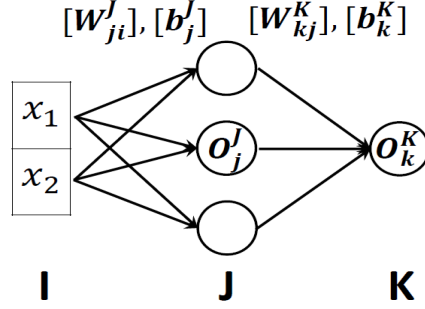


Figure 1: A Neural Network

ML Example: Consider a simple **Neural Network** (NN) as shown in Fig. 1 by Ryan Harris, *Neural Network Tutorial: The Back-Propagation Algorithm (2012)* on YouTube, where

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in R^2, \text{ input data point (vector)} \quad (11)$$

$$\mathbf{W}^J = [W_{ji}^J] = \begin{bmatrix} W_{11}^J & W_{12}^J \\ W_{21}^J & W_{22}^J \\ W_{31}^J & W_{32}^J \end{bmatrix} \in R^{3 \times 2}, \mathbf{b}^J = \begin{bmatrix} b_1^J \\ b_2^J \\ b_3^J \end{bmatrix} \in R^{3 \times 1} \quad (12)$$

$$\mathbf{W}^K = [W_{kj}^K] = [W_{11}^K \ W_{12}^K \ W_{13}^K] \in R^{1 \times 3}, \mathbf{b}^K = b_k^K \in R^1 \quad (13)$$

$$\mathbf{O}^J = \mathbf{W}^J \mathbf{x} = \begin{bmatrix} O_1^J \\ O_2^J \\ O_3^J \end{bmatrix} : \text{hidden layer } J, O_j^J: \text{hidden node} \quad (14)$$

$$y_k(\mathbf{W}) = O_k^K = \sigma(t) : \text{output value, } y_k \approx y_k^*: \text{true value} \quad (15)$$

$$\sigma(t) = \sigma(\mathbf{W}^K \mathbf{W}^J \mathbf{x}) = \frac{1}{1 + e^{-t}} : \text{sigmoid function, } t = \mathbf{W}^K \mathbf{O}^J \quad (16)$$

$$\sigma'(t) = \frac{d\sigma(t)}{dt} = \frac{e^{-t}}{(1 + e^{-t})^2} = \frac{1 + e^{-t} - 1}{(1 + e^{-t})^2} = \sigma(1 - \sigma). \quad (17)$$

$$\text{Minimize Error } E(\mathbf{W}) = \frac{1}{2} \sum_{k \in K} (y_k(\mathbf{W}) - y_k^*)^2, \forall \mathbf{W}. \quad (18)$$

$$\nabla E(\mathbf{W}) = \left\langle \dots, \frac{\partial E(\mathbf{W})}{\partial W_{ji}^J}, \dots, \frac{\partial E(\mathbf{W})}{\partial W_{kj}^K}, \dots \right\rangle \quad (19)$$

$$\frac{\partial E(\mathbf{W})}{\partial W_{ji}^J} = (y_k(\mathbf{W}) - y_k^*) \frac{\partial y_k(\mathbf{W})}{\partial W_{ji}^J}, \quad y_k(\mathbf{W}) = \sigma(\mathbf{W}) = \sigma(t) \quad (20)$$

$$\begin{aligned} \frac{\partial y_k(\mathbf{W})}{\partial W_{ji}^J} &= \frac{\partial \sigma(\mathbf{W}^K \mathbf{W}^J \mathbf{x})}{\partial W_{ji}^J} = \frac{d\sigma(t)}{dt} \frac{\partial (\mathbf{W}^K \mathbf{O}^J)}{\partial W_{ji}^J} \\ &= \sigma'(t) \frac{\partial}{\partial W_{ji}^J} \sum_{l \in J} W_{kl}^K O_l^J = \sigma' \sum_{l \in J} W_{kl}^K \frac{\partial O_l^J}{\partial W_{ji}^J} \\ &= \sigma' \sum_{l \in J} \left[W_{kl}^K \frac{\partial O_l^J}{\partial O_j^J} \frac{\partial}{\partial W_{ji}^J} \sum_{m \in I} W_{lm}^J x_m \right] = \sigma' W_{kj}^K x_i \end{aligned} \quad (21)$$

$$p_i = \frac{\partial E(\mathbf{W})}{\partial W_{ji}^J} = (y_k - y_k^*) O_k^K (1 - O_k^K) W_{kj}^K x_i, \quad \frac{\partial E(\mathbf{W})}{\partial b_j^J} = 1 \quad (22)$$

$$\frac{\partial y_k(\mathbf{W})}{\partial W_{kj}^K} = \sigma' \frac{\partial (\mathbf{W}^K \mathbf{O}^J)}{\partial W_{kj}^K} = \sigma' \frac{\partial}{\partial W_{kj}^K} \sum_{l \in J} W_{kl}^K O_l^J = \sigma' O_j^J \quad (22)$$

$$p_j = \frac{\partial E(\mathbf{W})}{\partial W_{kj}^K} = (y_k - y_k^*) O_k^K (1 - O_k^K) O_j^J, \quad \frac{\partial E(\mathbf{W})}{\partial b_k^K} = 1 \quad (23)$$

Back Propagate $(y_k - y_k^*) O_k^K (1 - O_k^K)$ from k Nodes to j Nodes!! (24)

$$\text{Method: } W_m^{(n)} = W_m^{(n-1)} - \alpha^{(n-1)} p_m^{(n-1)}, \quad m = i \text{ or } j, \forall m. \quad (25)$$

$$\mathbf{p} = \nabla E(\mathbf{W}) = \langle \dots, p_m, \dots \rangle, \quad \mathbf{W}^{(0)} \text{ and } \alpha^{(n-1)} \text{ given,}$$

$$\mathbf{W} = \langle \dots, W_m, \dots \rangle = \langle \dots, W_{ji}^J, \dots \rangle, \quad \Delta W_m = W_m^{(n)} - W_m^{(n-1)}.$$

Forward: $\nabla E(\mathbf{W}) = \left\langle \frac{\partial E(\mathbf{W})}{\partial W_1}, \dots, \frac{\partial E(\mathbf{W})}{\partial W_{1,000,000}} \right\rangle = \langle \dots, p_m, \dots \rangle$

\Rightarrow One Million Paths from j Nodes to k Nodes

\Rightarrow One Million Differentiations!!!

\Rightarrow Need Output Error $(y_k - y_k^*)$ One Million Times (Individually).

\Rightarrow Need 116 Days!!

Backward: One Output Error $(y_k - y_k^*)$ from k Nodes to j Nodes

\Rightarrow Get One Million $\frac{\partial E(\mathbf{W})}{\partial W_m}$ Layer-by-Layer Simultaneously!!

\Rightarrow Need Only 1 Second, 10^7 Times Faster!!