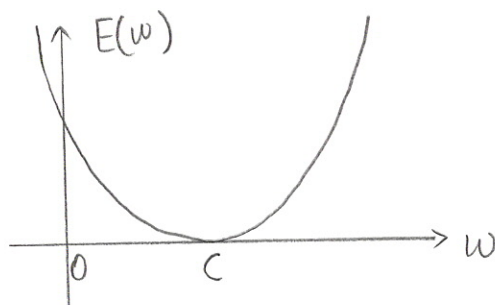


Gradient Descent

- Minimize an objective function.

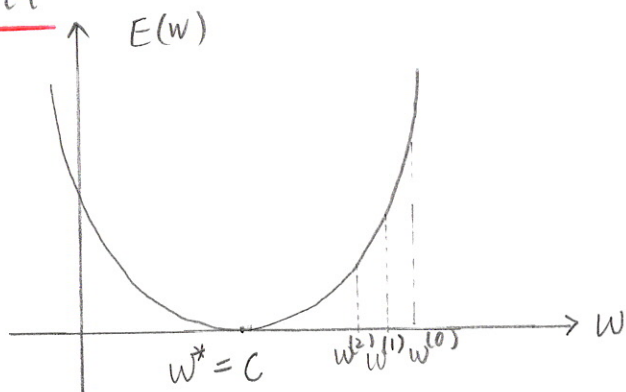
$$E(w) = \frac{1}{2} (w - c)^2$$



$$\frac{dE(w)}{dw} = \frac{1}{2} \cdot 2(w - c) = 0 \implies w^* = c. \quad \text{A closed-form solution}$$

For functions $E(w)$ that are difficult to calculate the closed-form solution, or when the closed-form solution does not exist, we use another method to find w^* .

Gradient Descent



$w^{(0)}$: an initial value for w

Update the new $w^{(1)}$ by

$$w^{(1)} = w^{(0)} + \eta \cdot p$$

p : a direction ; $\eta > 0$: a step size, or learning rate

It's a constant.

Assume we know $E(w^{(0)})$, we want to find $E(w^{(1)})$.

"Taylor Expansion" says:

$$E(w^{(1)}) \approx E(w^{(0)}) + \underbrace{\frac{dE(w)}{dw} \Big|_{w=w^{(0)}}}_{E'(w^{(0)})} \times (w^{(1)} - w^{(0)})$$

e.g. $E(w) = \frac{1}{2} (w - c)^2$

$$\frac{dE(w)}{dw} = \frac{1}{2} \cdot 2 \cdot (w - c) = w - c$$

$$\underbrace{\frac{dE(w)}{dw} \Big|_{w=w^{(0)}}}_{E'(w^{(0)})} = w^{(0)} - c$$

(3)

$$E(w^{(1)}) \approx E(w^{(0)}) + (w^{(0)} - c) \times (w^{(1)} - w^{(0)}) \leq E(w^{(0)})$$

$$(w^{(0)} - c) \times (w^{(1)} - w^{(0)}) \leq 0 \quad \dots (*)$$

$$\text{since } w^{(1)} = w^{(0)} + \eta \cdot P \quad \dots (1)$$

plug (1) into (*)

$$(w^{(0)} - c) \times (w^{(0)} + \eta \cdot P - w^{(0)}) = (w^{(0)} - c) \times \eta \cdot P \leq 0$$

$$\therefore \eta > 0$$

$$\therefore P \cdot (w^{(0)} - c) \leq 0 \quad \dots (2)$$

$$\text{Let } P = -(w^{(0)} - c)$$

$$\text{then: LHS of (2)} = -(w^{(0)} - c)^2 \leq 0$$

conclusion:

$$P = - \left. \frac{dE(w)}{dw} \right|_{w=w^{(0)}} = -E'(w^{(0)})$$

$$w^{(1)} = w^{(0)} - \eta \cdot \left. \frac{dE(w)}{dw} \right|_{w=w^{(0)}}$$

In the τ -th iteration, update $w^{(\tau)}$ by

$$w^{(\tau)} = w^{(\tau-1)} - \eta \cdot \left. \frac{dE(w)}{dw} \right|_{w=w^{(\tau-1)}}$$

$\tau = 1, 2, 3, \dots$

For example: $E(w) = \frac{1}{2}(w-4)^2$

Initialize $w^{(0)} = 3$, let $\eta = 1e-2$, $\epsilon = 1e-8$

For $\tau = 1, 2, 3, \dots$

$$w^{(\tau)} = w^{(\tau-1)} - \eta \cdot [w^{(\tau-1)} - 4]$$

$$\text{update } E(w^{(\tau)}) = \frac{1}{2}(w^{(\tau)} - 4)^2$$

check stopping criterion:

$$\text{if } |E(w^{(\tau)}) - E(w^{(\tau-1)})| < \epsilon$$

break

Return w^*

Code: gradient-descent-parabola.py

Summary: Scalar-version Gradient-Descent Algorithm

$$\text{Objective: } w^* = \arg \min_w E(w)$$

Initialize: $w^{(0)}$, $E(w^{(0)})$, $\eta > 0$, $\epsilon > 0$, max Iter.

For $\tau = 1, 2, 3, \dots, \text{max Iter.}$

① Update gradient (or derivative)

$$\nabla_w E(w) \Big|_{w=w^{(\tau-1)}} = \dots$$

② update weight

$$w^{(\tau)} = w^{(\tau-1)} - \eta \cdot \nabla_w E(w) \Big|_{w=w^{(\tau-1)}}$$

③ update objective function.

$$E(w^{(\tau)}) = \dots$$

④ Check stopping criterion.

$$\text{If } \left| E(w^{(\tau)}) - E(w^{(\tau-1)}) \right| < \epsilon.$$

break.

Return $w^* = w^{(\tau_{\text{stop}})}$

When $E(\vec{w})$ is a function of a vector $\vec{w} = [w_0, w_1, \dots, w_m]^T$

(6)

update \vec{w} by

$$\vec{w}^{(\tau)} = \vec{w}^{(\tau-1)} + \eta \cdot \vec{p}$$

$$\vec{p} = -\nabla_{\vec{w}} E(\vec{w}) \Big|_{\vec{w} = \vec{w}^{(\tau-1)}}$$

The gradient of the function $E(\cdot)$ with respect to \vec{w} is:

$$\nabla_{\vec{w}} E(\vec{w}) \triangleq \begin{bmatrix} \frac{\partial E(\vec{w})}{\partial w_0} \\ \frac{\partial E(\vec{w})}{\partial w_1} \\ \vdots \\ \frac{\partial E(\vec{w})}{\partial w_m} \end{bmatrix}, \text{ where } \frac{\partial E(\vec{w})}{\partial w_m} \text{ is called the partial derivative.}$$

$m = 0, 1, 2, \dots, M$

$$\nabla_{\vec{w}} E(\vec{w}) \Big|_{\vec{w} = \vec{w}^{(\tau-1)}} \triangleq \begin{bmatrix} \frac{\partial E(\vec{w})}{\partial w_0} \Big|_{\vec{w} = \vec{w}^{(\tau-1)}} \\ \frac{\partial E(\vec{w})}{\partial w_1} \Big|_{\vec{w} = \vec{w}^{(\tau-1)}} \\ \vdots \\ \frac{\partial E(\vec{w})}{\partial w_m} \Big|_{\vec{w} = \vec{w}^{(\tau-1)}} \end{bmatrix}$$

Example:

$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

(7)

$$E(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (w_0 + w_1 x_n - t_n)^2$$

$$\frac{\partial E(\vec{w})}{\partial w_0} = \frac{1}{2} \sum_{n=1}^N 2 (w_0 + w_1 x_n - t_n) = \sum_{n=1}^N (w_0 + w_1 x_n - t_n)$$

$$\frac{\partial E(\vec{w})}{\partial w_1} = \frac{1}{2} \sum_{n=1}^N 2 (w_0 + w_1 x_n - t_n) \cdot x_n = \sum_{n=1}^N (w_0 + w_1 x_n - t_n) \cdot x_n$$

$$\nabla_{\vec{w}} E(\vec{w}) = \begin{bmatrix} \frac{\partial E(\vec{w})}{\partial w_0} \\ \frac{\partial E(\vec{w})}{\partial w_1} \end{bmatrix}$$

$$\nabla_{\vec{w}} E(\vec{w}) \Big|_{\vec{w} = \vec{w}^{(l-1)}} = \begin{bmatrix} \sum_{n=1}^N (w_0^{(l-1)} + w_1^{(l-1)} \cdot x_n - t_n) \\ \sum_{n=1}^N (w_0^{(l-1)} + w_1^{(l-1)} \cdot x_n - t_n) \cdot x_n \\ N \cdot w_0^{(l-1)} + w_1^{(l-1)} \cdot \sum_{n=1}^N x_n - \sum_{n=1}^N t_n \\ w_0^{(l-1)} \cdot \sum_{n=1}^N x_n + w_1^{(l-1)} \cdot \sum_{n=1}^N x_n^2 - \sum_{n=1}^N t_n \cdot x_n \end{bmatrix}$$

Summary : vector-version Gradient Descent Algorithm

$$\text{Objective : } \vec{w}^* = \arg \min_{\vec{w}} E(\vec{w})$$

Initialize : $\vec{w}^{(0)}$, $E(\vec{w}^{(0)})$, $\eta > 0$, $\epsilon > 0$, maxIter

For $\tau = 1, 2, 3, \dots, \text{maxIter}$

① Update gradient

$$\left. \nabla_{\vec{w}} E(\vec{w}) \right|_{\vec{w} = \vec{w}^{(\tau-1)}} = \dots$$

② update weights

$$\vec{w}^{(\tau)} = \vec{w}^{(\tau-1)} - \eta \cdot \left. \nabla_{\vec{w}} E(\vec{w}) \right|_{\vec{w} = \vec{w}^{(\tau-1)}}$$

③ update objective function

$$E(\vec{w}^{(\tau)}) = \dots$$

④ Check stopping criterion:

$$\text{If } \left| E(\vec{w}^{(\tau)}) - E(\vec{w}^{(\tau-1)}) \right| < \epsilon$$

break

Return $\vec{w}^* = \vec{w}^{(\tau_{\text{stop}})}$

Vector-version gradient descent for Linear Regression.

9

$$E(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (\vec{x}_n^T \cdot \vec{w} - t_n)^2$$

$$\vec{w}^{(l)} = \vec{w}^{(l-1)} - \eta \cdot \nabla_{\vec{w}} E(\vec{w}) \Big|_{\vec{w} = \vec{w}^{(l-1)}}$$

Recall: $\frac{d \vec{a}^T \cdot \vec{w}}{d \vec{w}} = \vec{a}$

Hence, $\nabla_{\vec{w}} E(\vec{w}) = \frac{1}{2} \cdot \sum_{n=1}^N 2 \cdot (\vec{x}_n^T \cdot \vec{w} - t_n) \cdot \frac{d \vec{x}_n^T \cdot \vec{w}}{d \vec{w}}$

$$= \sum_{n=1}^N (\vec{x}_n^T \cdot \vec{w} - t_n) \cdot \vec{x}_n = \sum_{n=1}^N \vec{x}_n \cdot (\vec{x}_n^T \cdot \vec{w} - t_n)$$

$$\nabla_{\vec{w}} E(\vec{w}) = \left(\sum_{n=1}^N \vec{x}_n \cdot \vec{x}_n^T \right) \vec{w} - \sum_{n=1}^N \vec{x}_n \cdot t_n$$

Recall: $X^T = [\vec{x}_1^T \ \vec{x}_2^T \ \dots \ \vec{x}_N^T]$, $X = \begin{bmatrix} \vec{x}_1^T \\ \vec{x}_2^T \\ \vdots \\ \vec{x}_N^T \end{bmatrix}$, $\vec{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}$

Hence, $\nabla_{\vec{w}} E(\vec{w}) = X^T \cdot X \cdot \vec{w} - X^T \cdot \vec{t}$

$$\nabla_{\vec{w}} E(\vec{w}) \Big|_{\vec{w} = \vec{w}^{(l-1)}} = X^T \cdot X \cdot \vec{w}^{(l-1)} - X^T \cdot \vec{t}$$

$$\Rightarrow \vec{w}^{(l)} = \vec{w}^{(l-1)} - \eta \cdot [X^T \cdot X \cdot \vec{w}^{(l-1)} - X^T \cdot \vec{t}]$$

$l = 1, 2, 3, \dots$