



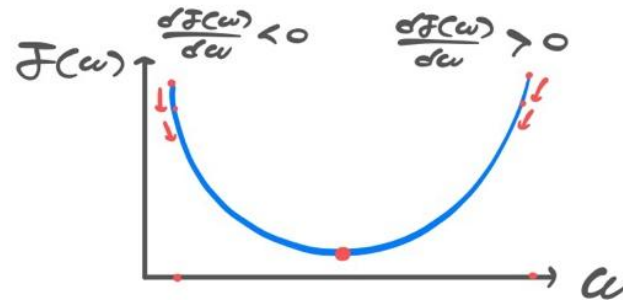
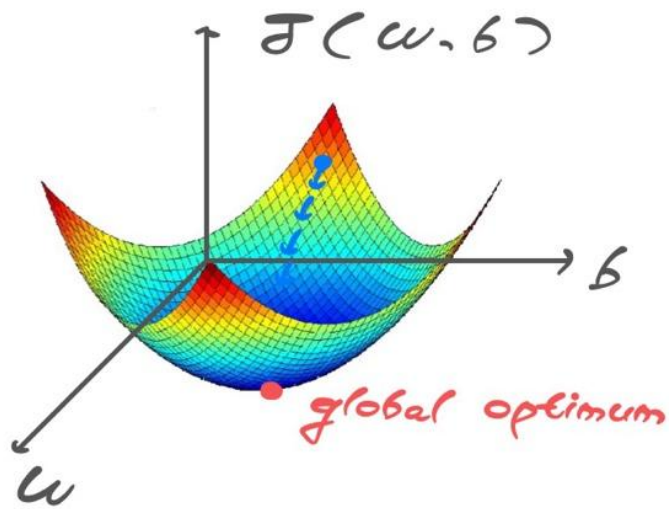
Gradient Descent

Convex Optimization

How to minimize cost function?

Gradient Descent

Want to find W, b that minimize $J(W, b)$



$$\begin{aligned} w &:= w - \alpha \frac{\partial J(w, b)}{\partial w} \\ b &:= b - \alpha \frac{\partial J(w, b)}{\partial b} \end{aligned}$$

update

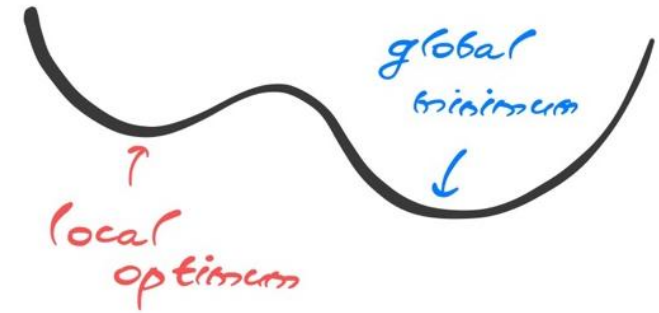
α : learning rate

Introduction

Gradient Descent

- Simple and gives some kind of meaningful result for both convex and nonconvex optimization
- > It tries to improve the function value by moving in a direction related to the gradient (i.e first derivative)

1. For convex optimization, it gives the global optimum
2. For nonconvex optimization, it arrives at a local optimum



Convex function

Function f is convex

ie " $f(dx + (1-d)y) \leq df(x) + (1-d)f(y)$ "

for $\forall x, y$ and $d \in [0, 1]$

Introduction

Taylor expansion of a function all of whose derivatives exist at x

$$f(x+h) = f(x) + hf'(x) + \frac{h^2}{2} f''(x) + \frac{h^3}{3!} f'''(x) \dots$$

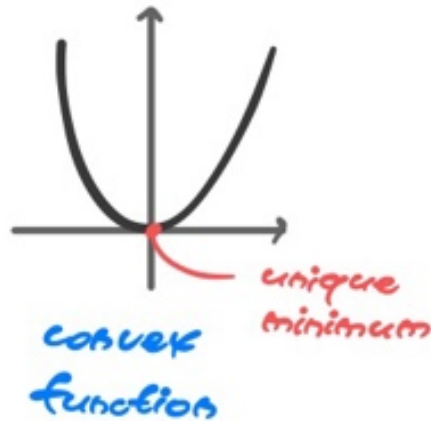
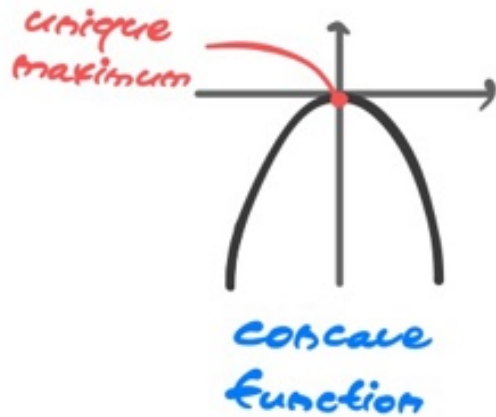
The function is convex

if $f''(x) \geq 0$ for $\forall x$

($\Leftrightarrow f'(x)$ is increasing function of x)

Introduction

The minimum is attained when $f'(x) = 0$
since $f'(x)$ keeps increasing to the left and right of that
thus the global minimum is unique.



Introduction

Gradient Descent

- Given a differentiable function of $f(x)$ and an initial parameter of x_1
- Iteratively moving the parameter to the lower value of $f(x)$
- By taking the direction of the negative gradient of $f(x)$

Why this works?

$$f(x) = f(a) + \frac{f'(a)}{1!} (x-a) + O(\|x-a\|^2)$$

Assume $a = x_i$ and $x = x_i + h u_i$ unit direction vector for the partial deriv

$$f(x_i + h u_i) = f(x_i) + h f'(x_i) u_i + h^2 O(1)$$

$$f(x_i + h u_i) - f(x_i) \approx h f'(x_i) u_i$$

Introduction

$$\begin{aligned} u^* &= \operatorname{argmin}_u \{f(x_i + hu) - f(x_i)\} \\ &= \operatorname{argmin}_u h f'(x_i) u = - \frac{f'(x_i)}{|f'(x_i)|} \end{aligned}$$

$$(\because f(x_i + hu) \leq f(x_i), \quad a \cdot b = |a||b| \cos \alpha)$$

$$x_{t+1} \leftarrow x_t + h u^* = x_t - h \frac{f'(x_t)}{|f'(x_t)|}$$

Introduction

Example

$$f(x_1, x_2) = (1-x_1)^2 + 100(x_2-x_1^2)^2$$

$$\frac{\partial}{\partial x_1} f(x_1, x_2) = -2(1-x_1) - 400x_1(x_2-x_1^2)$$

$$\frac{\partial}{\partial x_2} f(x_1, x_2) = 200(x_2-x_1^2) \quad \text{global minimum} = 0 \\ \text{at } (1, 1)$$

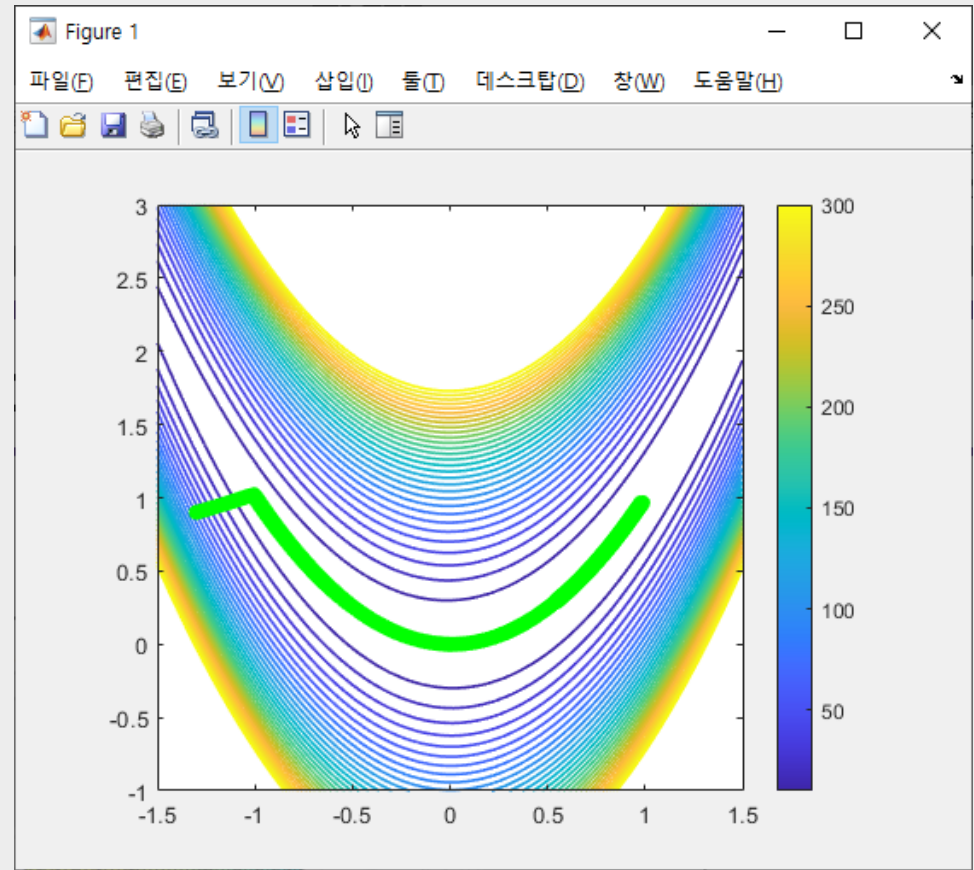
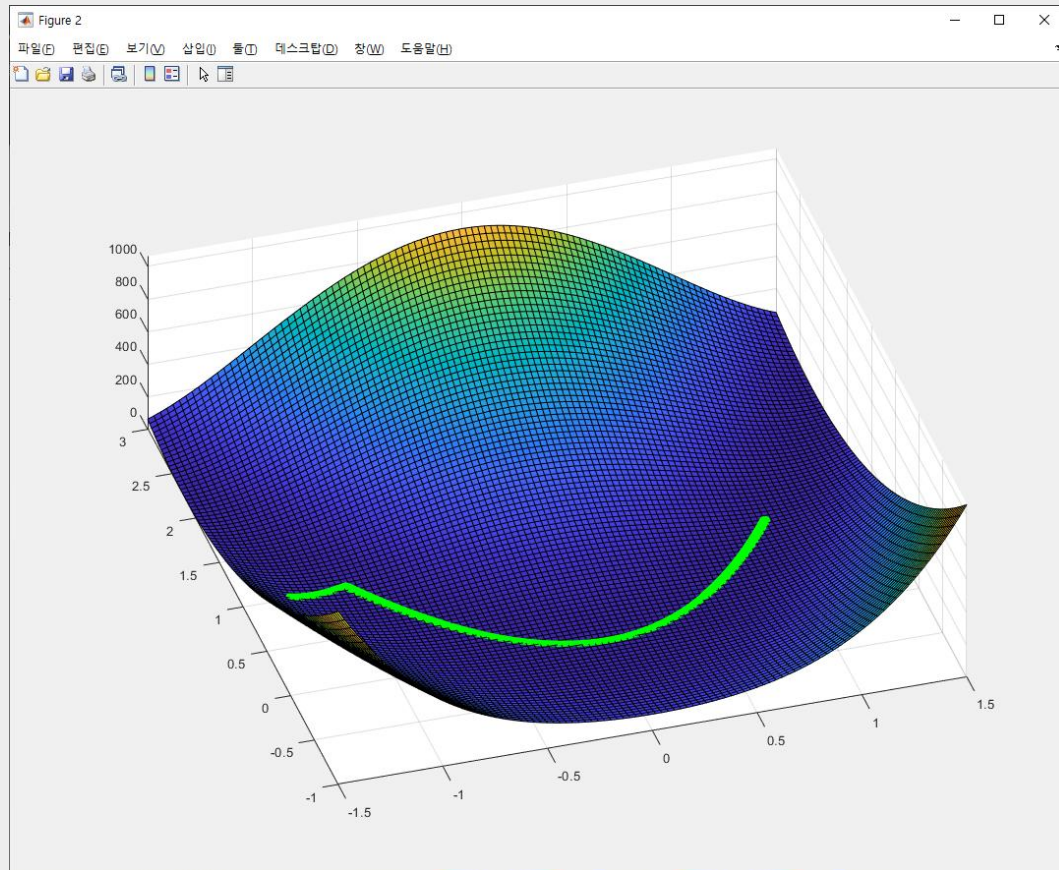
$$\text{initial point } \mathbb{x}^0 = (x_1^0, x_2^0) = (-1.3, 0.9)$$

$$f'(\mathbb{x}^0) = \left(\frac{\partial}{\partial x_1} f(x_1, x_2), \frac{\partial}{\partial x_2} f(x_1, x_2) \right) \\ = (-415.4, -158)$$

$$\mathbb{x}' \leftarrow \mathbb{x}^0 - \frac{f'(\mathbb{x}^0)}{|f'(\mathbb{x}^0)|}$$

repeat the update until converge

Introduction



1. Gradient descent for convex functions : univariate case

We start at some x_0 and
at step i update x_i to $x_{i+1} = x_i + \eta f'(x_i)$ for some $\eta < 0$
i.e move in direction where f decreases.

$$f(x_{i+1}) = f(x_i) + \eta f'(x_i) + \frac{\eta^2}{2} f''(x_i) \quad \text{ignore terms that involve } \eta^3 \text{ or higher}$$

Newton's method

Best value of η most reduction in one step

" $\eta = - \frac{f'(x)}{f''(x)}$ " which gives $f(x_{i+1}) = f(x_i) - \frac{(f'(x_i))^2}{2f''(x_i)}$

→ the algorithm makes progress so long as $f''(x_i) > 0$

2. Convex multivariate functions

If higher derivatives also exist, the multivariate Taylor expansion for an n -variate function f is

$$f(x+y) = f(x) + \nabla f(x) \cdot y + y^T \nabla^2 f(x) y + \dots$$

Hessian

$$H(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

- f is convex if the Hessian is positive semidefinite (i.e. $y^T \nabla^2 f y \geq 0$ for $\forall y$)

3. Gradient Descent for Constrained Optimization

- Constrained optimization consists of solving the following where K is a convex set and f is convex function.

$$\text{" } \min f(x) \text{ s.t. } x \in K \text{ "}$$

Lagrange Multiplier

new
optimization
problem

$$\min f(x) \text{ s.t. } g(x) = c$$

$$\exists d \in \mathbb{R} \text{ s.t. } \min f(x) + d(g(x) - c)$$

References

- [1] Boyd, Stephen, Stephen P. Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [2] 인공지능 및 기계학습 개론 I, https://www.edwith.org/machinelearning1_17/joinLectures/9738