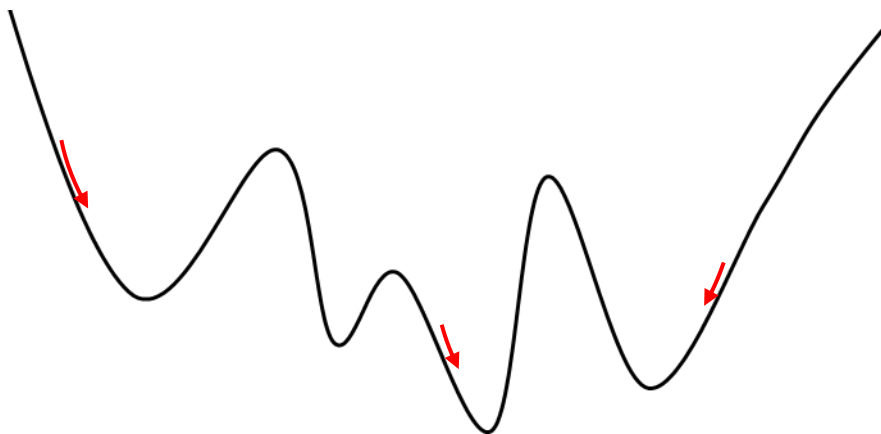


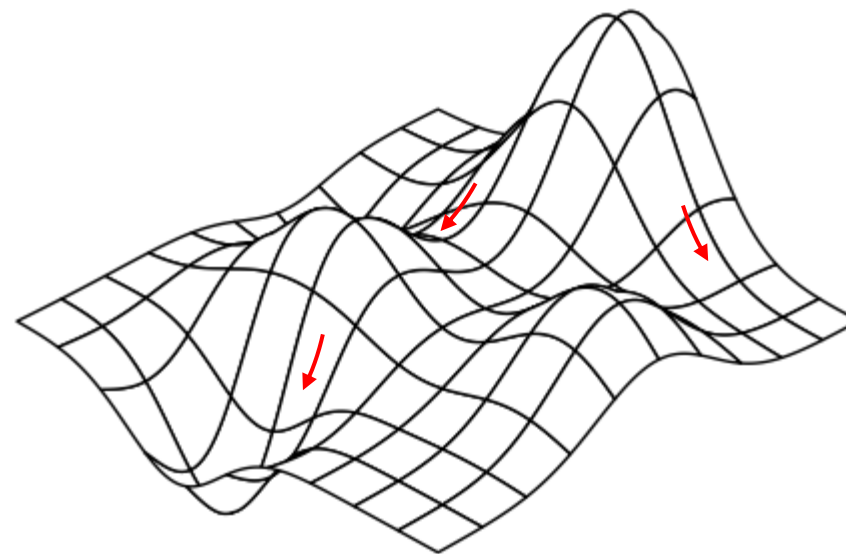
Chapter 03. 쉽게 배우는 경사 하강 학습법

STEP2. 심화 경사 하강 학습법

비볼록 함수 Non-convex Function



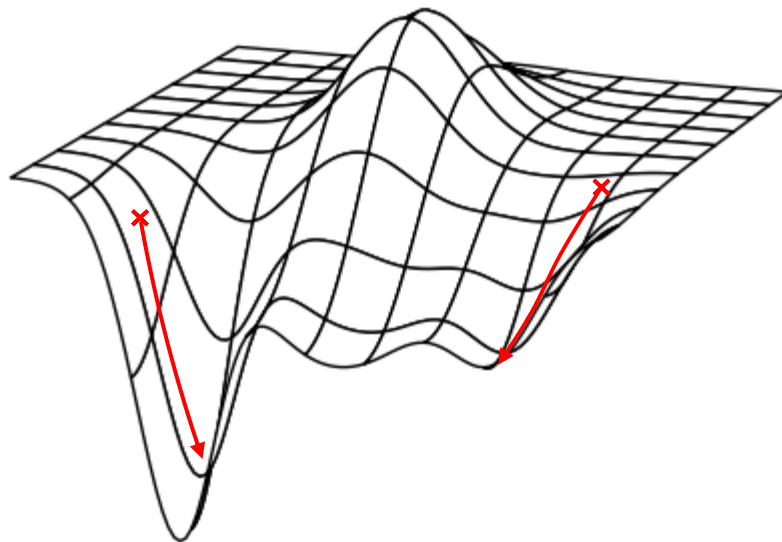
1-D Non-convex function



2-D Non-convex function

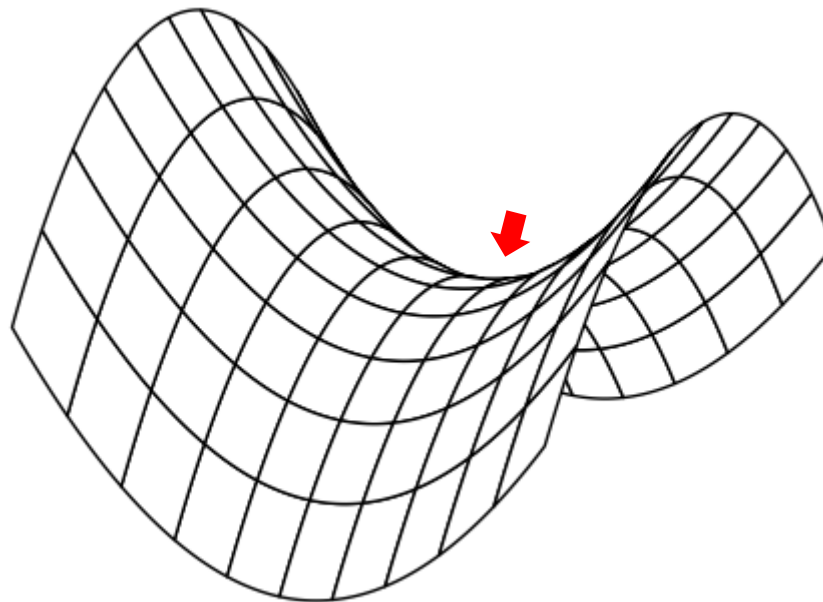
우리가 마주칠 대부분의 문제는 **비볼록 함수**이므로, 단순한 경사 하강법으로는 한계가 있다.

지역 최솟값 Local Minimum



경사 하강법을 사용할 경우, 초기값에 따라 Local minimum에 빠질 위험이 있다.

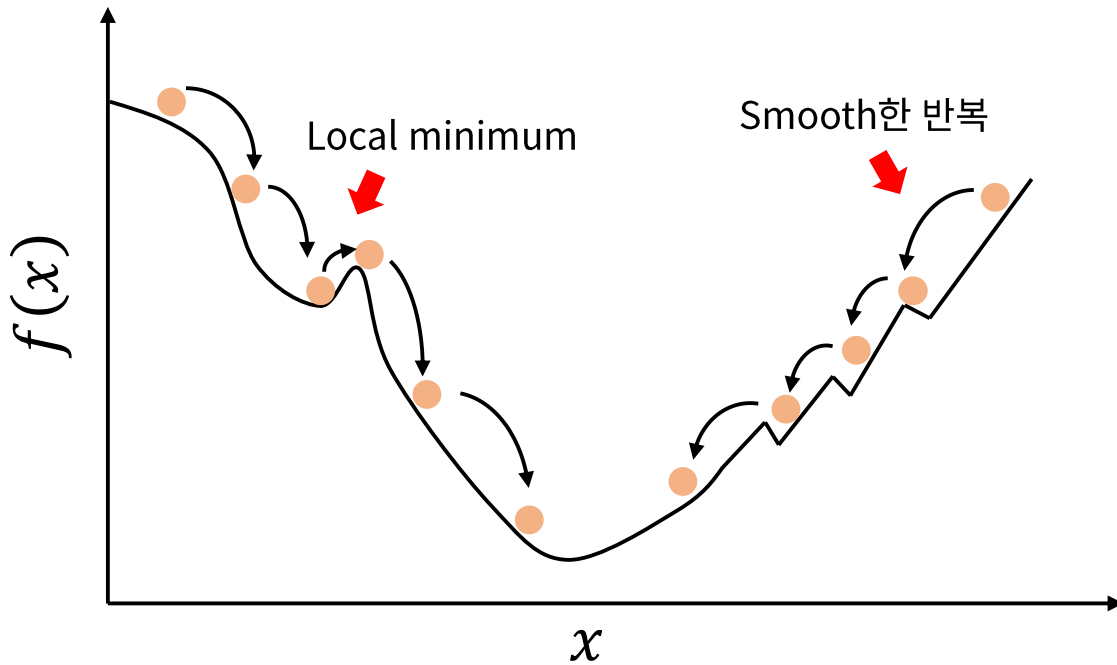
안장점 Saddle Point



안장점(Saddle point)은 기울기가 0이 되지만 극값이 아닌 지점을 말한다.
경사 하강법은 안장점에서 벗어나지 못한다.

관성 Momentum

돌이 굴러 떨어지듯, 이동 벡터를 이용해 이전 기울기에 영향을 받도록 하는 방법



$$\mathbf{v}_t = \gamma \mathbf{v}_{t-1} + \eta \nabla f(\mathbf{x}_{t-1})$$

$$\mathbf{x}_t = \mathbf{x}_{t-1} - \mathbf{v}_t$$

γ : 관성 계수 (momentum term) ≈ 0.9

η : 학습율 (Learning rate)

\mathbf{v}_t : t번째 step에서의 \mathbf{x} 의 이동 벡터

관성(Momentum)을 이용하면 **Local minimum과 잡음에 대처**할 수 있다.

이동 벡터(\mathbf{v}_t)를 추가로 사용하므로, 경사 하강법 대비 2배의 메모리를 사용한다.

적응적 기울기 AdaGrad

적응적 기울기 (Adaptive gradient; AdaGrad) : 변수별로 학습율이 달라지게 조절하는 알고리즘

$$\begin{array}{c}
 \text{Element-wise 제곱} \\
 \downarrow \\
 \boxed{
 \begin{array}{l}
 \mathbf{g}_t = \mathbf{g}_{t-1} + (\nabla f(\mathbf{x}_{t-1}))^2 \\
 \mathbf{x}_t = \mathbf{x}_{t-1} - \frac{\eta}{\sqrt{\mathbf{g}_t + \epsilon}} \cdot \nabla f(\mathbf{x}_{t-1})
 \end{array}
 } \\
 \uparrow \\
 \text{Element-wise 곱}
 \end{array}$$

\mathbf{g}_t : t번째 step까지의 기울기의 누적 크기

ϵ : 0으로 나누는 것을 방지하는 작은 값 $\approx 10^{-6}$

기울기가 커서 **학습이 많이 된 변수는 학습율을 감소**시켜, 다른 변수들이 잘 학습되도록 한다.

\mathbf{g}_t 가 계속해서 커져서 학습이 오래 진행되면 **더이상 학습이 이루어지지 않는 단점**이 있다.

RMSProp

RMSProp : AdaGrad의 문제점을 개선한 방법으로, 합 대신 지수평균을 사용

$$\begin{aligned} \mathbf{g}_t &= \gamma \mathbf{g}_{t-1} + (1 - \gamma)(\nabla f(\mathbf{x}_{t-1}))^2 \\ \mathbf{x}_t &= \mathbf{x}_{t-1} - \frac{\eta}{\sqrt{\mathbf{g}_t + \epsilon}} \cdot \nabla f(\mathbf{x}_{t-1}) \end{aligned}$$

γ : 지수 평균의 업데이트 계수

변수 간의 상대적인 학습율 차이는 유지하면서 \mathbf{g}_t 가 무한정 커지지 않아 학습을 오래 할 수 있다.

Adam

Adaptive moment estimation (Adam) : RMSProp과 Momentum의 장점을 결합한 알고리즘

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla f(\mathbf{x}_{t-1})$$

$$\mathbf{g}_t = \beta_2 \mathbf{g}_{t-1} + (1 - \beta_2) (\nabla f(\mathbf{x}_{t-1}))^2$$

$$\hat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \beta_1^t}, \hat{\mathbf{g}}_t = \frac{\mathbf{g}_t}{1 - \beta_2^t}$$

$$\mathbf{x}_t = \mathbf{x}_{t-1} - \frac{\eta}{\sqrt{\hat{\mathbf{g}}_t + \epsilon}} \cdot \hat{\mathbf{m}}_t$$

$$\beta_1 \approx 0.9$$

$$\beta_2 \approx 0.999$$

$$\epsilon \approx 10^{-8}$$

$\hat{\mathbf{m}}, \hat{\mathbf{g}}$: 초기값이 0인 것을 고려하여 보정한 값

Adam 최적화 방법은 **가장 최신의 기술(state-of-the-art)**이며, 딥러닝에서 가장 많이 사용된다.