

Gradient Descent Algorithm

and few more optimization...

Prof. Je-Won Kang
Electronic & Electrical Engineering
Ewha Womans University

Overview

- We have some function (loss function)

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

and want

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1) \leftarrow \text{object 함수 최소화 하는 } \theta_0, \theta_1 \text{ 찾는.}$$

Algorithm outline:

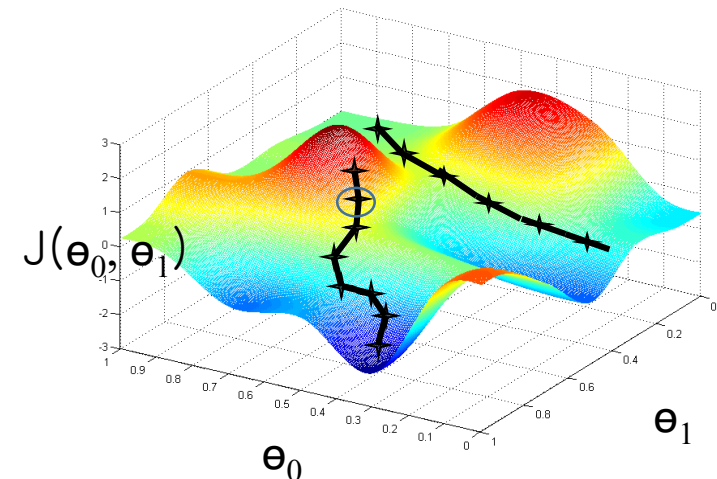
- Start with some initial parameters θ_0, θ_1
- Keep changing the parameter to reduce the loss function until we hopefully end up at a minimum.

Gradient Descent Algorithm

Key components

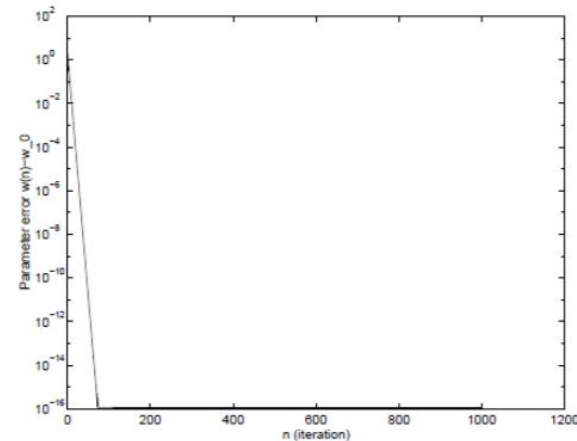
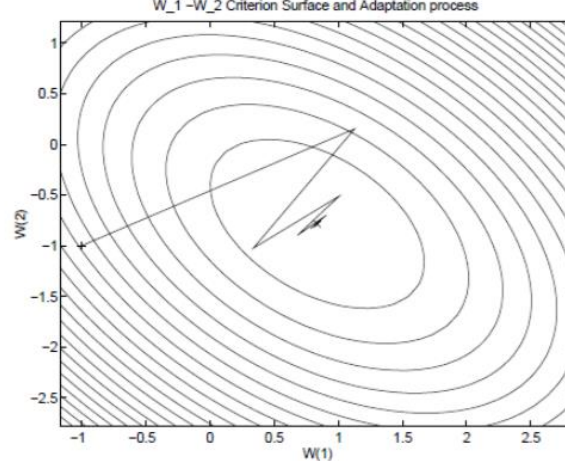
- **Gradient**: the derivative of vector functions (partial derivative along each dimension)
 - Direction of greatest increase (or decrease) of a function
- **The step size α** affects the rate at which the weight vector moves down the error surface and must be a positive number. (**hyper parameter**)
- θ is the **learnable parameters**
- The function J is **the objective function** that we want to minimize.

$$\text{repeat until convergence } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \end{array} \right. \}$$

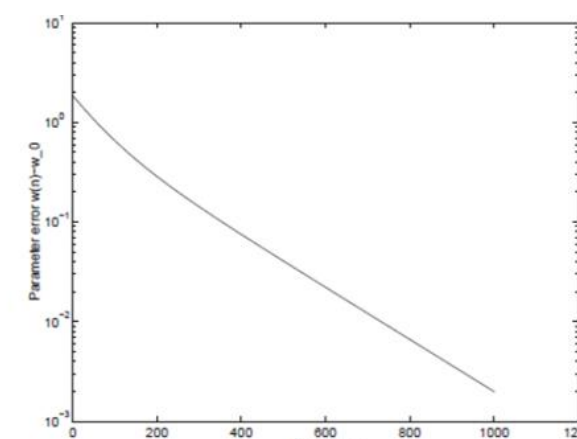
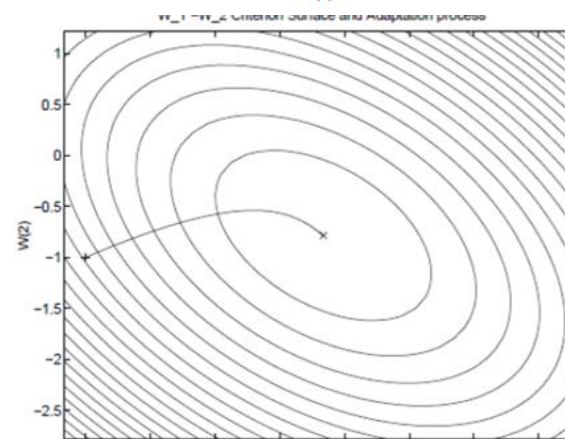


- If α is too small, gradient descent can be slow.
- If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

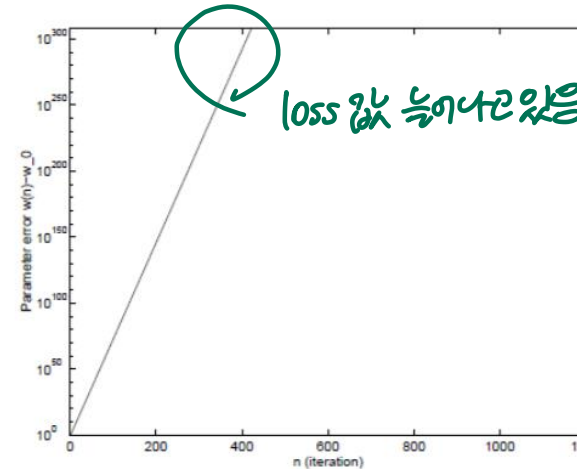
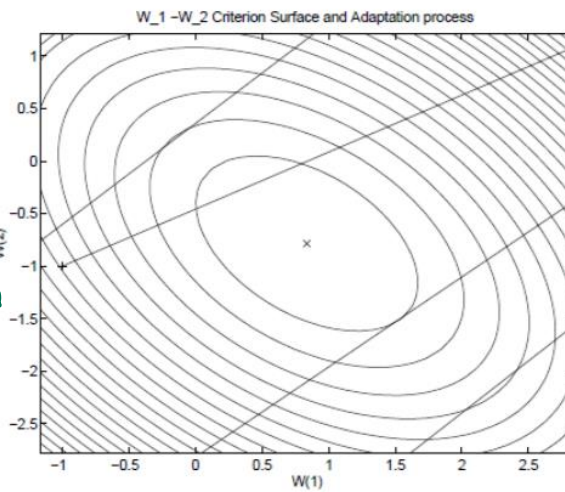
공간잡 →



α 작을 때 →
수렴속도 느림
but 안정적



α 큼 →
최소점 찾기 어렵



Batch gradient descent

θ_0 및 θ_1 업데이트
전체 샘플 m 개 모두 고려

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

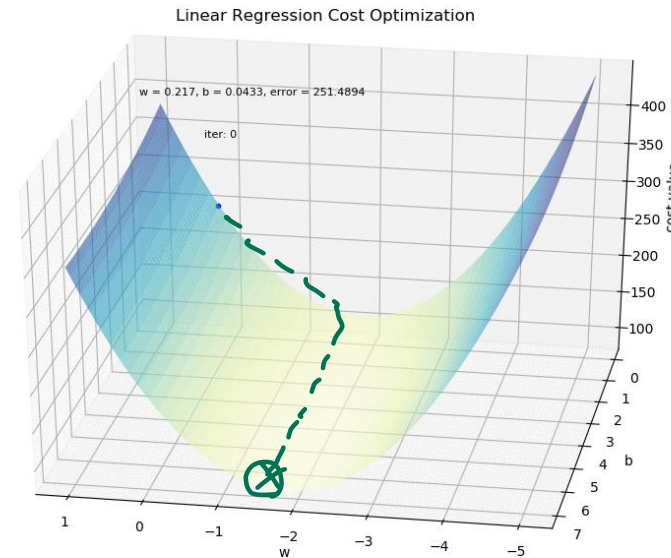
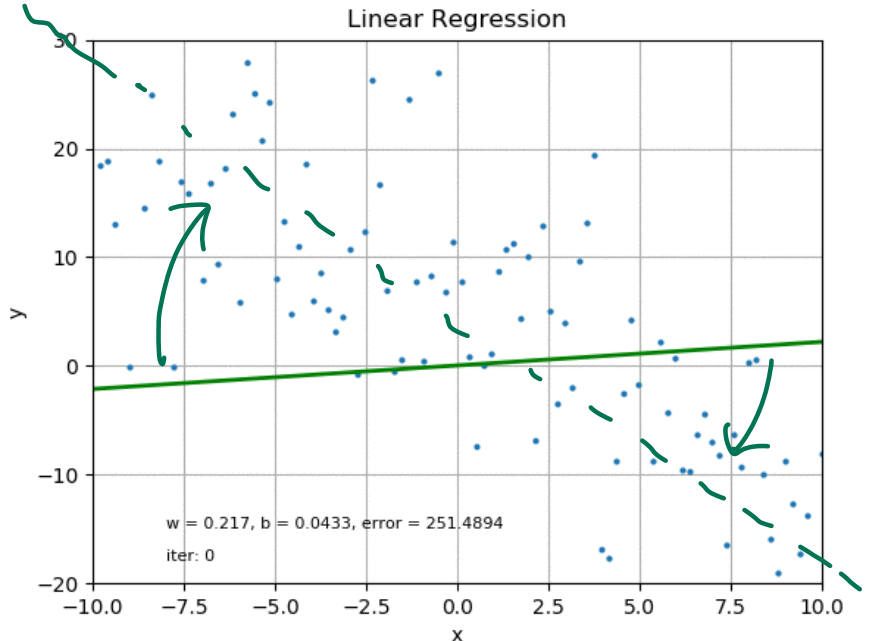
}

Linear regression model

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

What if the number of sample sizes m is increasing?



↓ 해결, but noise 영향 받기 쉬운 단점
↓ 방법의

Stochastic gradient descent (SGD)

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

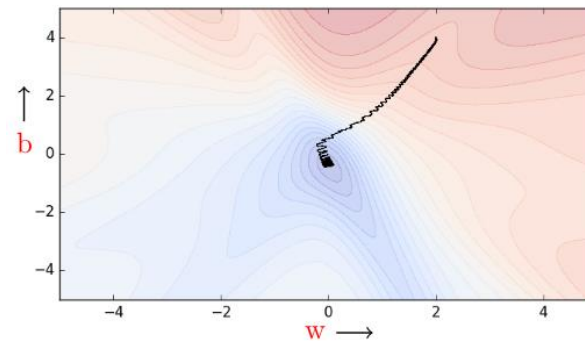
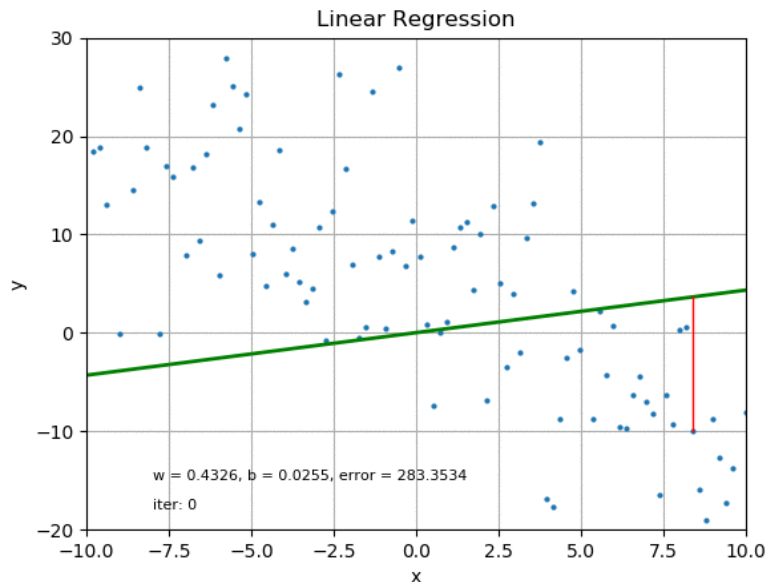
}

SGD : $m=1$

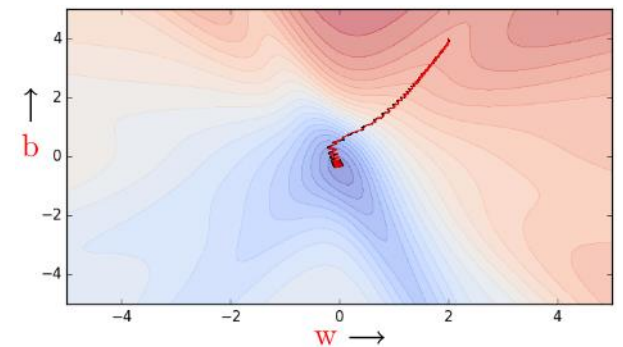
$1 < m < N$: mini-batch SGD

VS SGD : alleviate randomness

VS GD : less time in converging



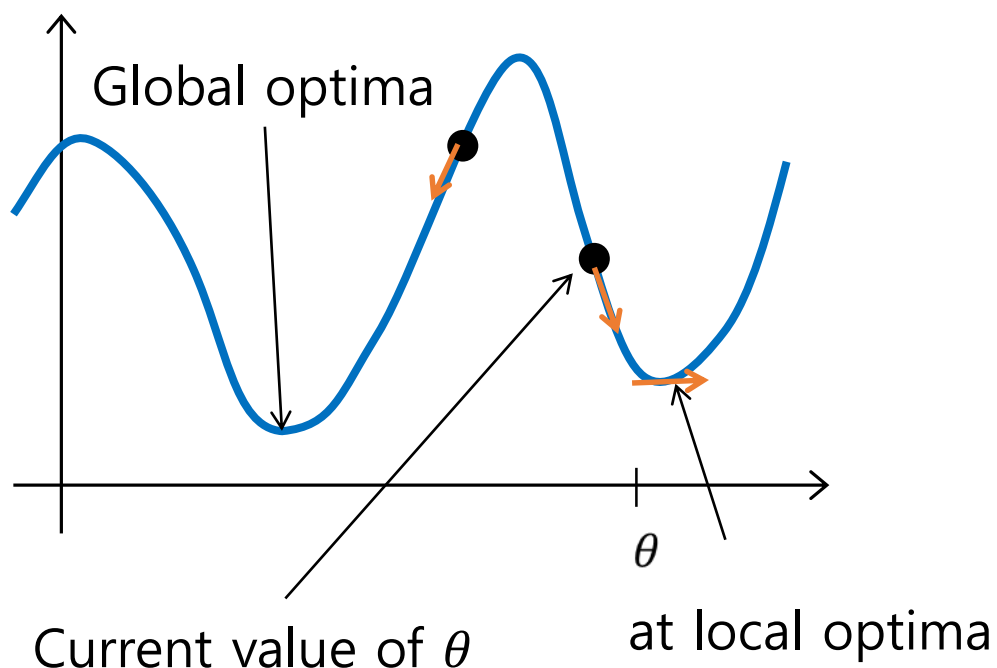
Many oscillations with greedy decision making



Less oscillations can be seen for $m=2$

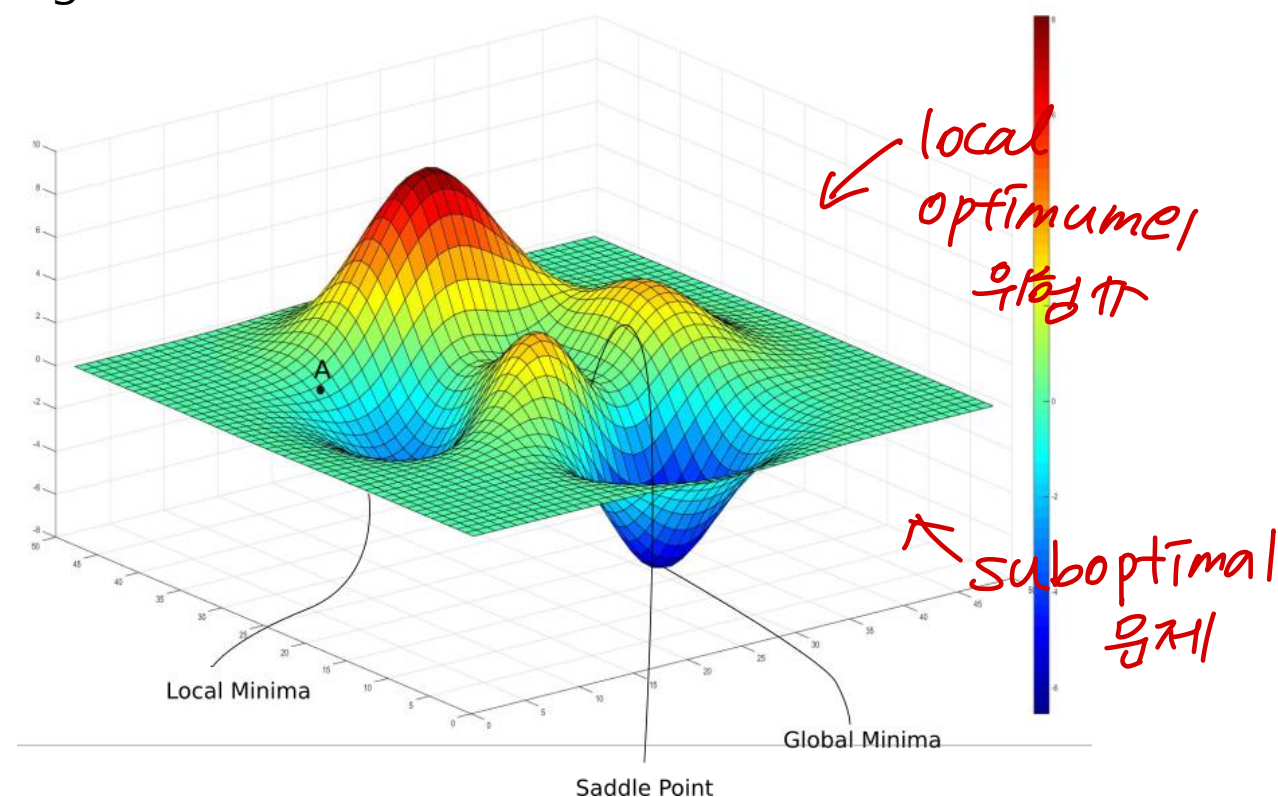
Limitation : Local Optimum

$$\theta_{new} \leftarrow \theta_{old} - \alpha \frac{\partial}{\partial \theta} J(\theta)$$



Cannot guarantee global minimum but attempt to find a good local minimum

Critical points with zero slope : $\nabla_x f(x)=0$
gives no information about which direction to move



Some ideas to avoid local minimum

Method of momentum

- SGD : very popular but tends to be slow and difficult to reach the minimum
- Method of momentum
 - Designed to speed up learning in high curvature and small/noise gradients
 - Exponentially weighted moving average of past gradients (low pass filtering)

← suboptimal 문제 해결 방법.

$$v_t = \begin{cases} g_1, & t = 1 \\ \rho v_{t-1} + (1 - \rho) g_t, & t > 1 \end{cases}$$

v_t : Exponentially weighted moving average at time t
 g_t : observation gradient at time t
 ρ (0~1): degree of weighting decrease (smoothing factor)

c.f $v_t = \rho^k v_{t-k} + (1 - \rho)[g_t + \rho g_{t-1} + \dots + \rho^{k-1} g_{t-k+1}]$

현재 과거 ← 먼 과거의 gradient는 더 작아짐.

SGD

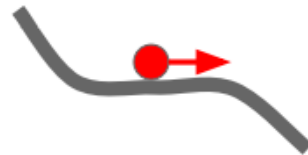
$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} J(\theta_t)$$

$$\text{let } g = \nabla_{\theta} J(\theta)$$

Local Minima



Saddle points



SGD + momentum :

← 과거값 반영해주는것.

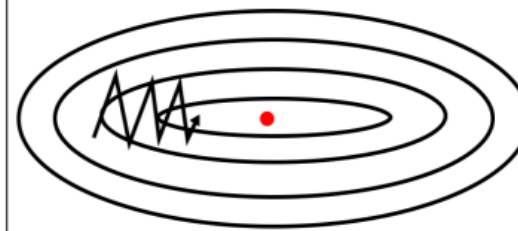
Use a velocity as a weighted moving average of previous gradients

$$v \leftarrow \rho v - \alpha g$$

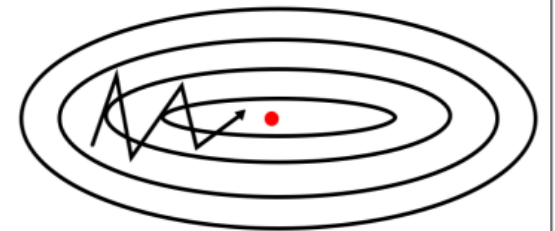
gradient 값이 0 이더라도 학습계속 진행할수 있음.

$$\theta \leftarrow \theta + v$$

SGD without momentum



SGD with momentum



A parameter is updated by linear combination of gradient and velocity

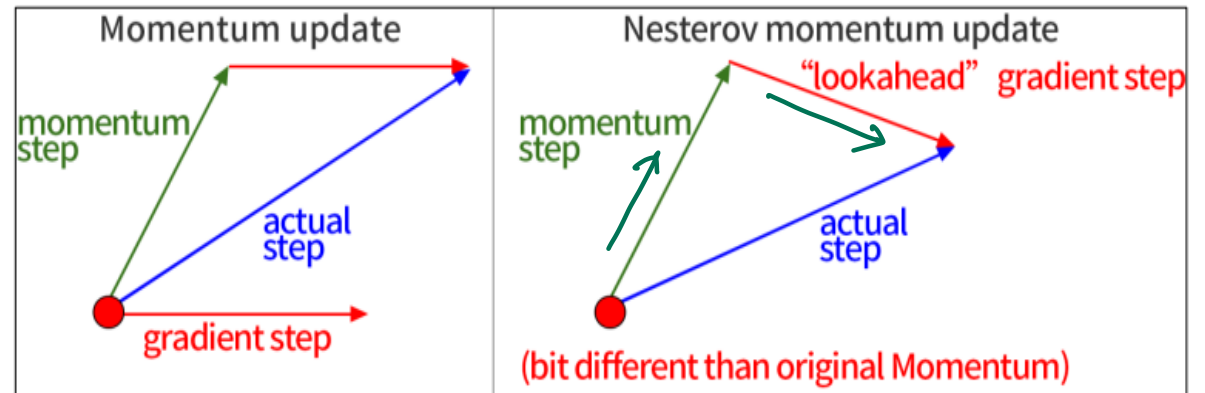
Some ideas to avoid local minimum

Method of momentum

- **Nesterov Momentum** ↖ 더 newest ⇒ gradient 먼저 평가하고 update 함
 - Difference from standard momentum: where gradient g is evaluated (i.e. "lookahead" gradient step)

$$v \leftarrow \rho v_t - \alpha \nabla_{\theta} J(\theta + \rho v)$$

$$\theta \leftarrow \theta + v$$



Some ideas to avoid local minimum per-parameter adaptive learning rates

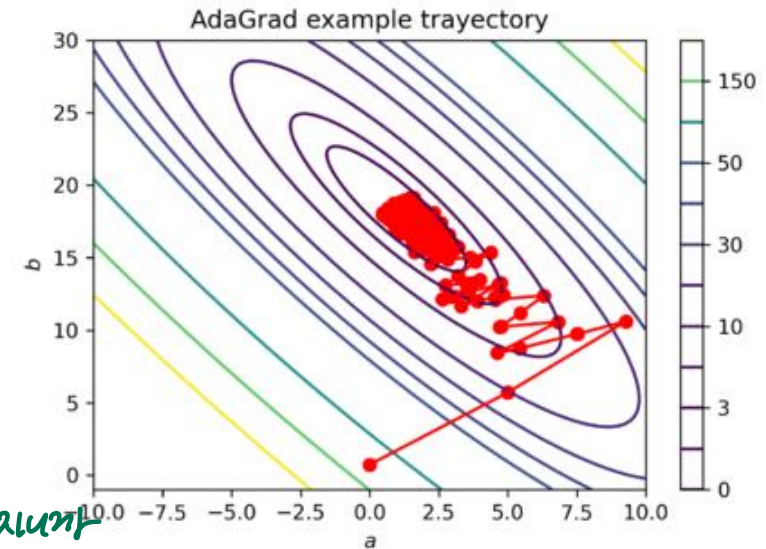
- **AdaGrad**: Adapts an individual learning rate of each direction
 - Slow down the learning rate when an accumulated gradient is large
 - Speed up the learning rate when an accumulated gradient is small
- Allows an automatic tuning of the learning rate per parameter

$$\text{let } g = \nabla_{\theta} J(\theta)$$

각 방향으로
learning rate를
적응적으로 조절하여
학습 효율 높이는 것

$$\begin{aligned} r &\leftarrow r + g \cdot g \\ \Delta \theta &\leftarrow \frac{\alpha}{\sqrt{r + \epsilon}} \cdot g \\ \theta &\leftarrow \theta - \Delta \theta \end{aligned}$$

gradient · gradient 값이 누적
→ r 점점 커짐
→ r 점점 작아짐
learning rate
조절 됨. → 이거 문제임. learning rate 계속 작아지니까
학습이 일어나지 않게 됨



Some ideas to avoid local minimum

per-parameter adaptive learning rates

- **RMSProp** : attempts to fix the drawbacks of AdaGrad, in which the learning rate becomes infinitesimally small and the algorithm is no longer able learning when the accumulated gradient is large.
- **(Remedy)** Gradient accumulation by weighted decaying

AdaGrad

$$r \leftarrow r + g \cdot g$$

$$\Delta \theta \leftarrow \frac{\alpha}{\sqrt{r + \varepsilon}} \cdot g$$

$$\theta \leftarrow \theta - \Delta \theta$$

RMSProp

$$r \leftarrow \rho r + (1 - \rho) g \cdot g$$

$$\Delta \theta \leftarrow \frac{\alpha}{\sqrt{r + \varepsilon}} \cdot g$$

$$\theta \leftarrow \theta - \Delta \theta$$

AdaGrad 문제 해결

극단적으로
학습속도 줄어드는거
막아줌.

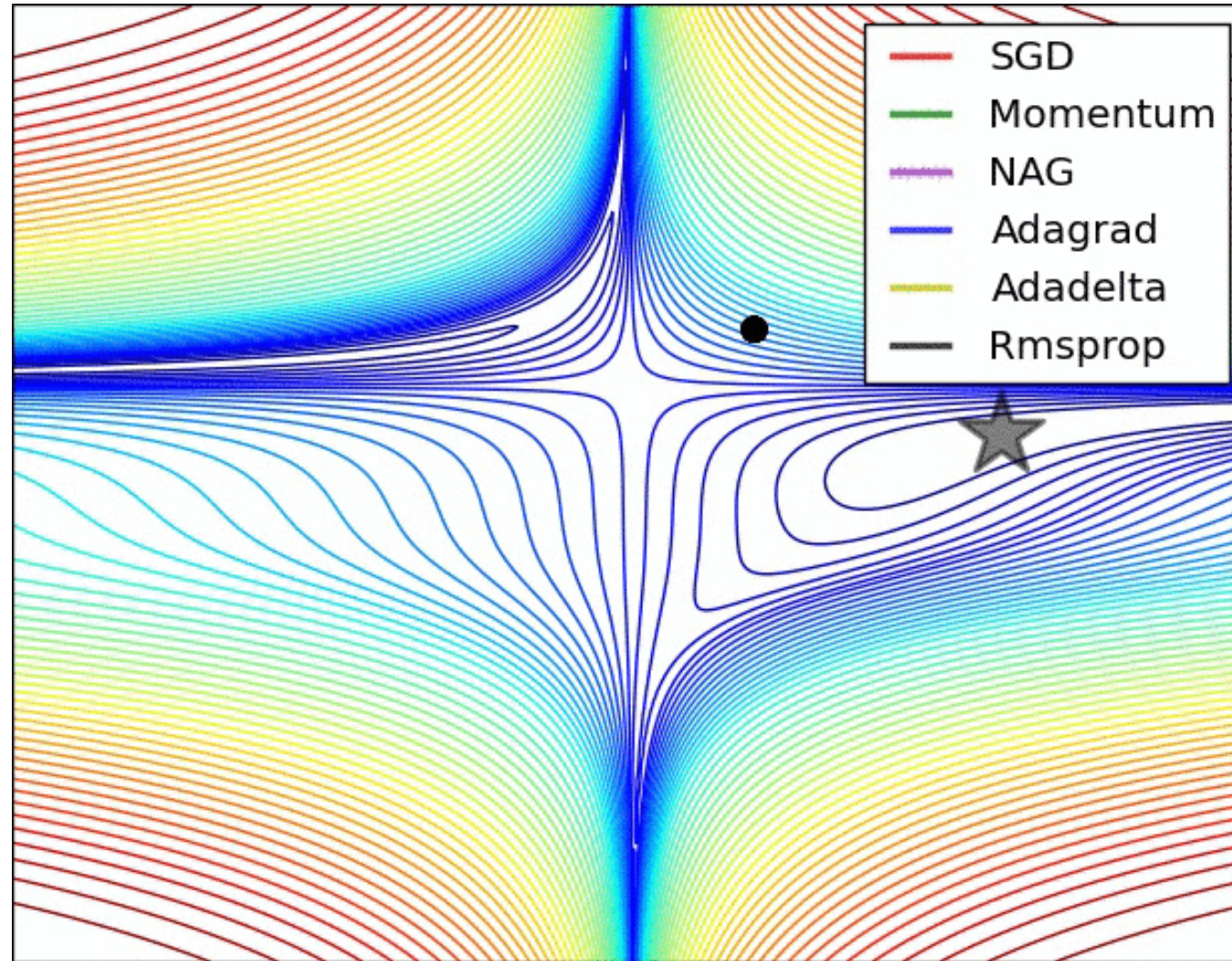
Some ideas to avoid local minimum

per-parameter adaptive learning rates

- ★ Adam (adaptive moment estimation) : RMSProp + momentum

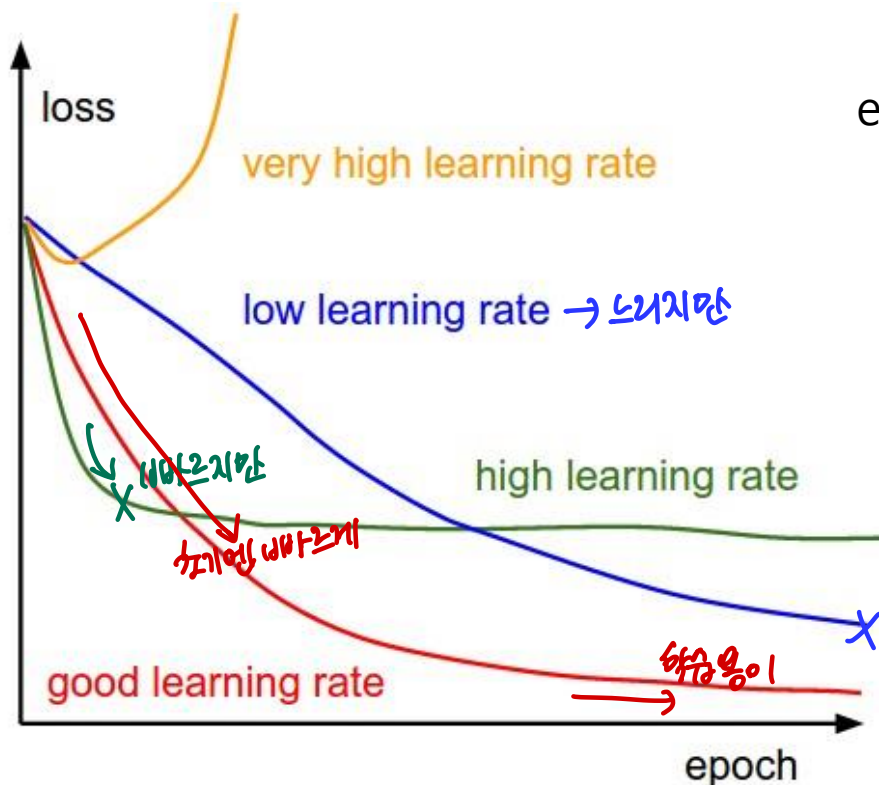
$$\text{let } g = \nabla_{\theta} J(\theta)$$

- ① Compute the first moment from momentum $s \leftarrow \rho_1 s + (1 - \rho_1) g \cdot g$
- ② Compute the second moment from RMSProp $r \leftarrow \rho_2 r + (1 - \rho_2) g \cdot g$
- ③ Correct the bias $s' \leftarrow \frac{s}{1 - \rho_1} \quad r' \leftarrow \frac{r}{1 - \rho_2}$
- ④ Update the parameters $\theta \leftarrow \theta - \alpha \frac{s'}{\sqrt{r' + \varepsilon}}$



Learning rate scheduling

- Learning rate: key hyper parameter for gradient-based algorithms ; need to gradually decrease learning rate over time



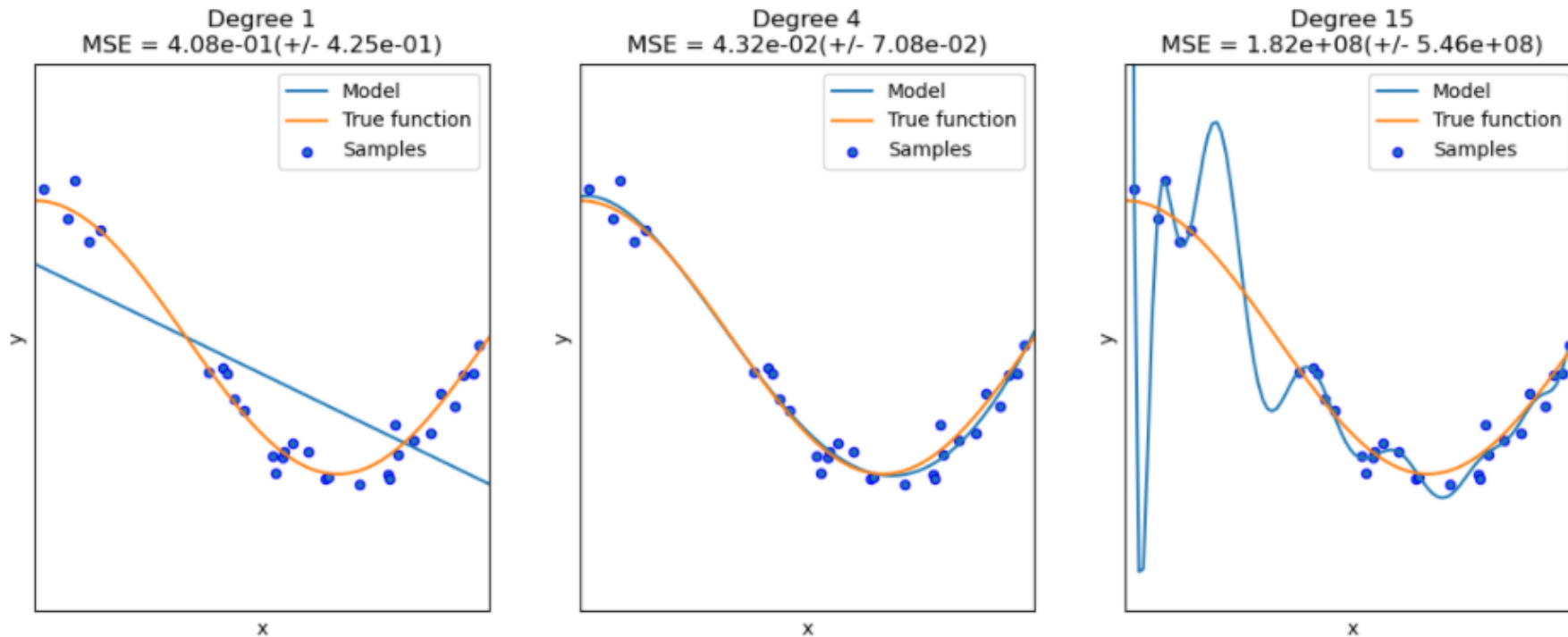
e.g. decay learning rate by every half of few epochs

$$\alpha = \alpha_0 e^{-kt}$$

$$\alpha = \frac{\alpha_0}{(1 + kt)}$$

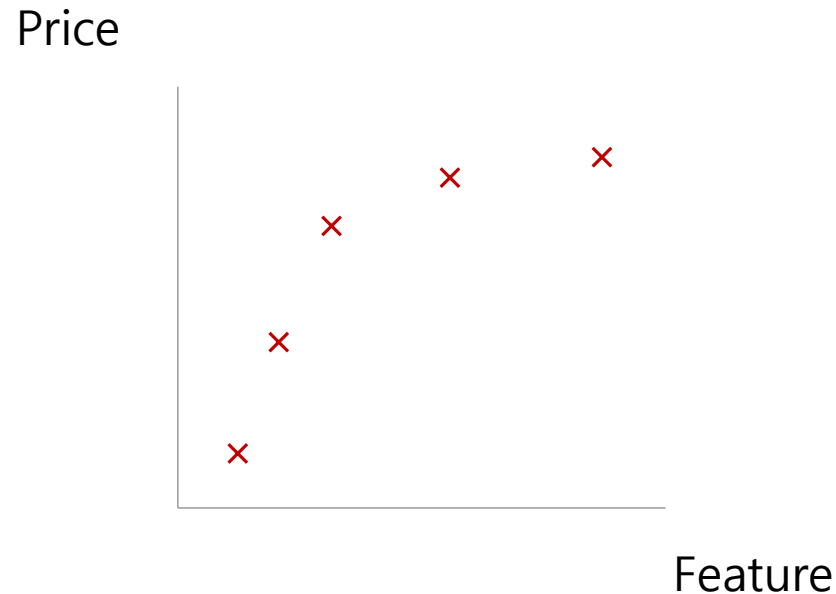
Some optimization in regression to avoid overfitting

- If we have too many features, the hypothesis may fit the training set very well. However, it may fail to generalize to new samples



Some optimization in regression to avoid overfitting

x_1 = size of house
 x_2 = no. of bedrooms
 x_3 = no. of floors
 x_4 = age of house
 x_5 =
 x_6 = kitchen size
 \vdots
 x_{100}



More features \rightarrow more parameters \rightarrow need more data ; (in practice) less data \rightarrow overfitting
Furthermore, mean-squared-error is sensitive to outliers

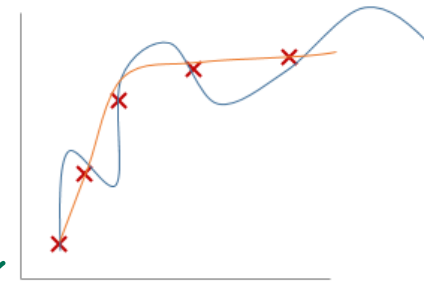
Some optimization in regression to avoid overfitting

1. Reduce number of features.
 - select which features to keep.
2. **Regularization.** ← 복잡한 모델 사용하게 되도 학습과정에서 모델의 복잡도에 대한 penalty를 줘서 모델이 overfitting 되지 않도록 해줌.
 - keep the features but reduce magnitude/values of parameters
 - Simple hypothesis and less prone to overfitting and robust to noise

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m \underbrace{(h_{\theta}(x^{(i)}) - y^{(i)})^2}_{\text{mse. "fitting"}} + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$\min_{\theta} J(\theta)$ λ : regularization parameter

← θ_j 가 클수록
늘어나는 오류
모델은 최대한
 θ_j 를 0으로 만들기
parameter
작게 만든다



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Quiz

What answers are correct? Select all that apply.

A. Gradient descent produces a numerical solution but it may not achieve a global optimum

Correct. The solution can be different with an initial point of an error surface

B. Momentum of previous gradient descent can help avoid overfitting

False. Momentum can avoid local minimum and help obtain a solution close to global optimum

C. Regularization can penalize the importance of some input features to avoid overfitting

Correct. Regularization can decrease some weights to have compact sets of parameters

Summary

- Optimization in general ML/DL
 - General ideas of gradient descent algorithm
 - Mostly stochastic gradient descent and its variants using gradient estimates
 - Adam is a good default choice in most cases
- Regularization
 - Reduce magnitude/values of parameters while keeping features
 - Simple hypothesis and less prone to overfitting
 - Robust to noise

Reference

- Book: Pattern Recognition and Machine Learning (by Christopher M. Bishop)
- Book: Machine Learning: a Probabilistic Perspective (by Kevin P. Murphy)
- <https://www.andrewng.org/courses/>