# Advanced Gradient Descent Methods
## Regularization

# CONTENTS

# CONTENTS

Batch
Gradient
Descent

**BGD**

# BGD

Based on Batch(total Data)

$$W := W - \alpha \frac{1}{m} \sum_{i=1}^{m} (Wx^i - y^i) x^i$$

# Stochastic Gradient Descent

**SGD**

# SGD

Based on (One Data)

$$W := W - \alpha \frac{1}{m} \left( W x^i - y^i \right) x^i$$
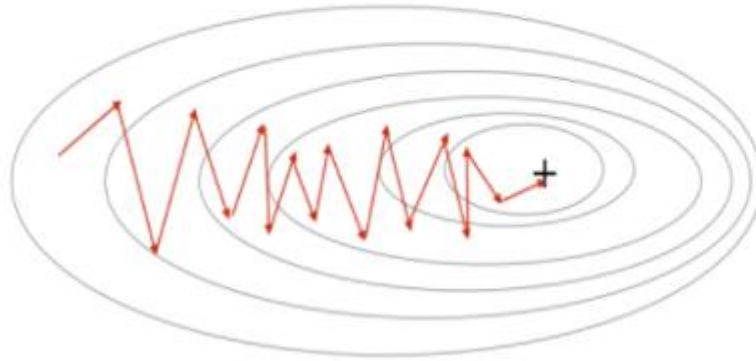
Mini-Batch
Gradient
Descent

**MSGD**

# MSGD
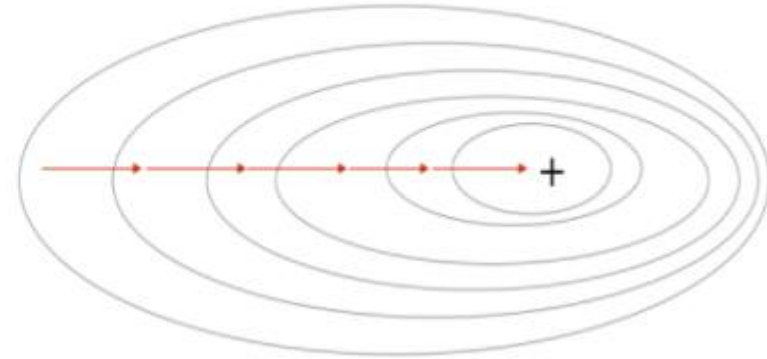
Based on Mini Batch (One Data)

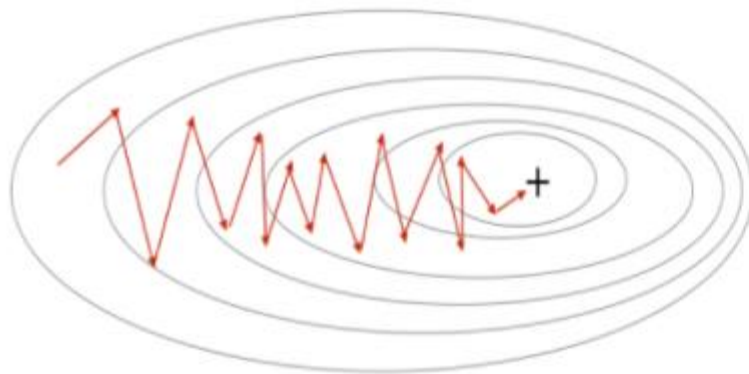$$W := W - \alpha \frac{1}{m} \sum_{i=1}^{m} (Wx^i - y^i) x^i$$
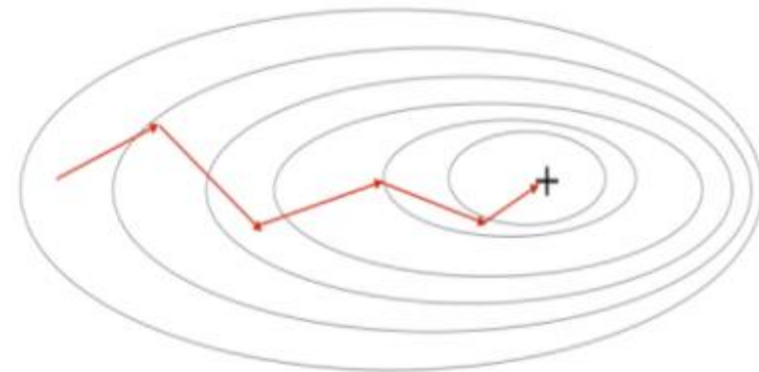
# Visualize

Stochastic Gradient Descent

Gradient Descent

Stochastic Gradient Descent

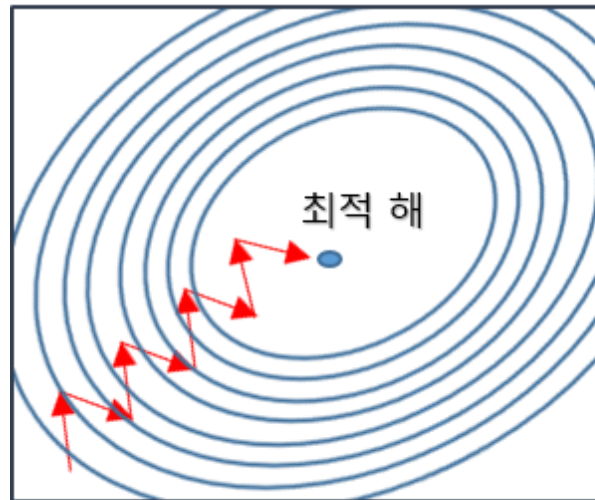Mini-Batch Gradient Descent

# Momentum

# Momentum

Like Acceleration => Rolling a ball

$$V(t) = m * V(t-1) - \alpha \frac{\partial}{\partial w} Cost(w) \quad\text{————————}\quad V(0)=0$$

$$W(t+1)=W(t)+V(t) \longleftarrow \quad W=weight$$

# Momentum



확률적 경사 하강법

모멘텀

Nesterov
Accelrated
Gradient

**NAG**

# NAG

# In Momentum Step, gradient



Difference between Momentum and NAG. Picture from CS231.

# NAG

$$V(t)=m*V(t-1)-\alpha \frac{\partial}{\partial (w+m*V(t-1))}Cost(w)$$

$$W(t+1)=W(t)+V(t)$$

NAG

**Adaptive Gradient**

**Adagrad**

# Adagrad

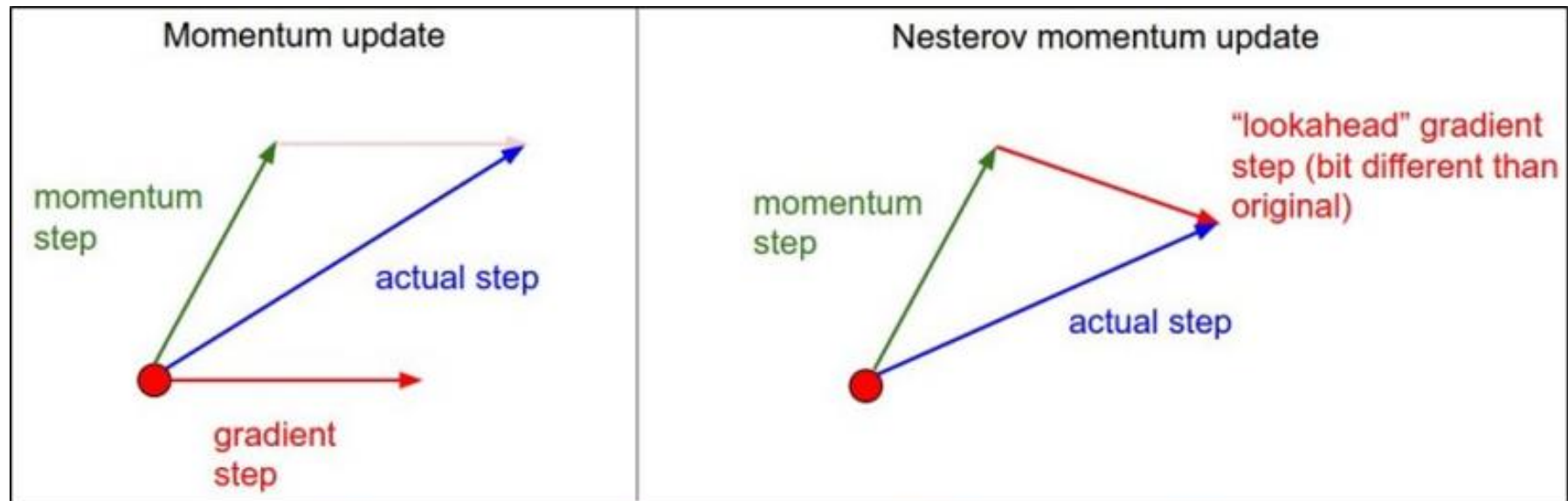$$G(t)=G(t-1)+\left(\frac{\partial}{\partial w(t)}Cost(w(t))\right)^2 \quad G(t) = \text{Vector W[i] element}$$

$$W(t+1)=W(t)-\alpha\frac{1}{\sqrt{G(t)+\epsilon}}\frac{\partial}{\partial w(i)}Cost(w(i)) \quad W(t) = \text{Vector W[i] element}$$

# Adagrad

# Problem

1. $G(0)=0$ and $G(t)=0$, insert $\in$

2. Infinite Training $G(t)$ is infinite

# Adagrad

Adaptive
Gradient

RMSProp

# RMSProp

It Complements the adagrad

$$G(t) = \gamma\, G(t-1) + (1-\gamma)\left(\frac{\partial}{\partial w(t)} Cost(w(t))\right)^2 \quad \text{G(t) = Vector W[i] element}$$
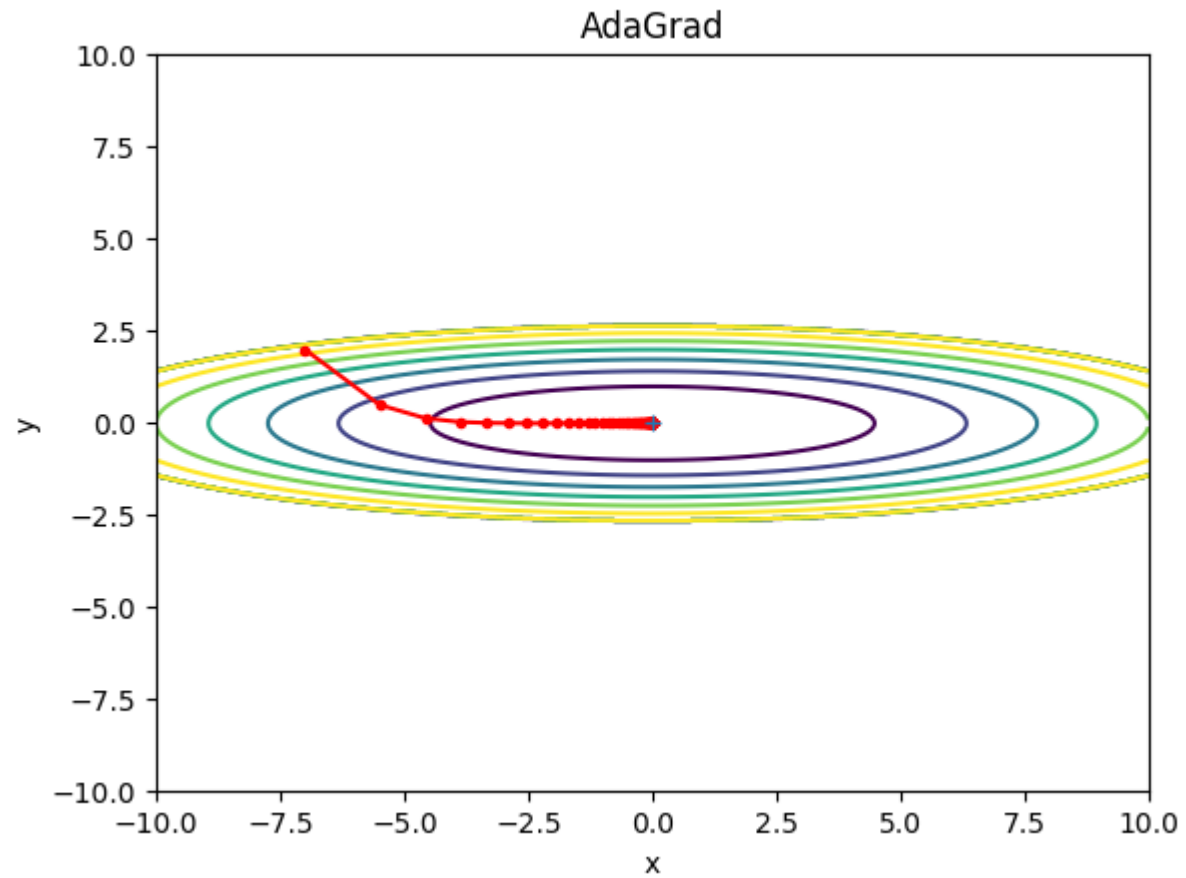
$$W(t+1) = W(t) - \alpha \frac{1}{\sqrt{G(t)+\epsilon}} \frac{\partial}{\partial w(i)} Cost(w(i)) \quad \text{W(t) = Vector W[i] element}$$

$$\gamma\text{`s Value} \Rightarrow 0.9\text{~}0.999$$

Adaptive
Gradient

AdaDelta

# Adadelta

Update Parameter W, W`s unit ?
Remove Learning Rate

$$G(t) = \gamma G(t-1) + (1-\gamma)\left(\frac{\partial}{\partial w(t)} Cost(w(t))\right)^2 \quad G(t) = \text{Vector W[i] element}$$

$$S(t) = \gamma S(t-1) + (1-\gamma)\Delta\theta^2$$

Hesian Matrix

$$W(t+1) = W(t) \frac{\sqrt{S(t)+\epsilon}}{\sqrt{G(t)+\epsilon}} \frac{\partial}{\partial w(i)} Cost(w(i)) \quad W(t) = \text{Vector W[i] element}$$

$\gamma$`s Value => 0.9~0.999

# Animation

애니메이션
http://shuuki4.github.io/deep%20learning/2016/05/20/Gradient-Descent-
Algorithm-Overview.html

# Adaptive Moment Estimation

**Adam**

# Adam

Moment + Adaptive

Moment is not Momentum Probability Moment

What is Moment ? => Kocw 김충락 교수님(수리통계학)

1-Moment=>E[X]

Not Known Moment => Estimation

2-Moment=>E[$X^2$]

# Adam

$$m_t = \beta_1 m_{t-1} + (1 + \beta_1)g_t$$

$$v_t = \beta_2 v_{t-1} + (1 + \beta_2)g_t^2$$

If Initial m, v is 0 , weight=>Zero biased

If decay rate is small, $(\beta_1, \beta_2 \ close \ one)$ weight=>biased
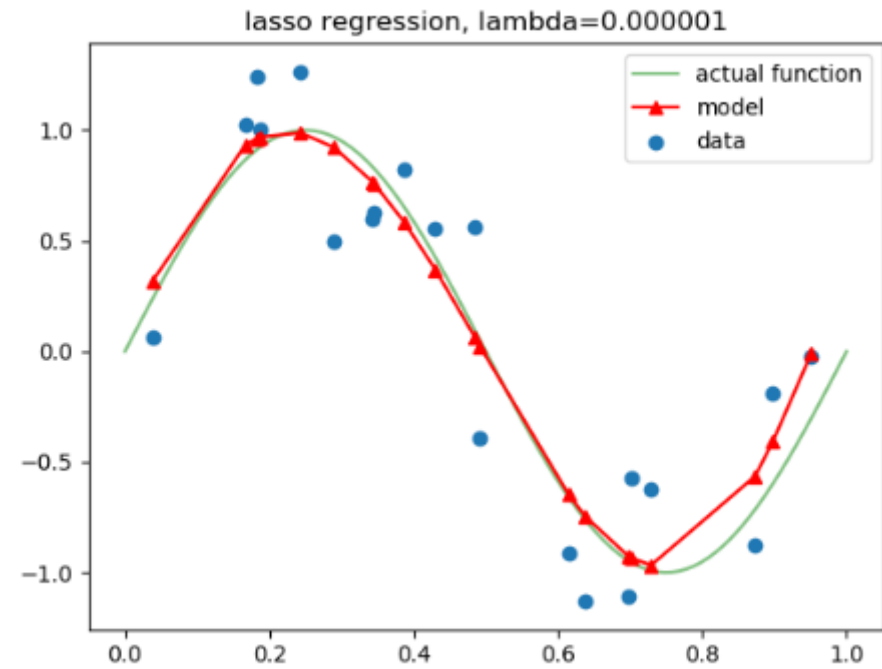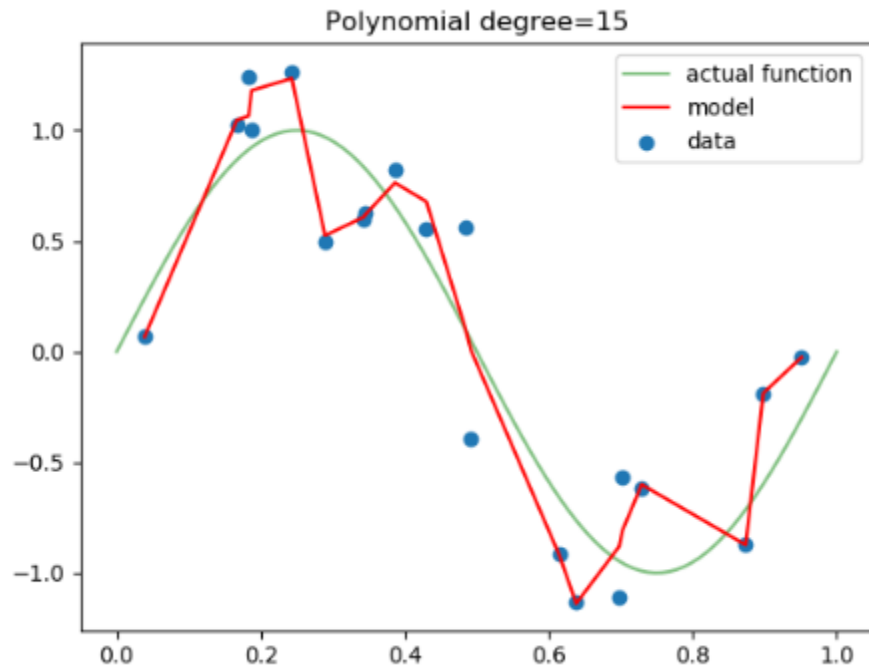
$$m_t = \beta_1 m_{t-1} + (1 + \beta_1)g_t$$

# Regularization

# Regularization

Lasso

**L1-Norm**

# L1−Norm

## Lasso

$$J(\theta) = \frac{1}{2m}\sum\left(h(x^i) - y^i\right)^2$$

$$J(\theta) = \frac{1}{2m}\sum\left(h(x^i) - y^i\right)^2 + \frac{\tau}{2}\sum\left|\theta_j\right|$$

Ridge

**L2-Norm**

# L2-Norm

## Ridge

$$J(\theta) = \frac{1}{2m}\sum\left(h(x^i) - y^i\right)^2$$

$$J(\theta) = \frac{1}{2m}\sum\left(h(x^i) - y^i\right)^2 + \frac{\tau}{2}\sum\theta_j^2$$

# L2-Norm

## Ridge

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

$$\theta_j := \theta_j (1 - \alpha \frac{\tau}{m}) - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^i) - y^i \right) x_j^i$$

감사합니다
THANK YOU