

[2-1. Optimization Algorithms]

< Mini-batch Gradient Descent >

배치 경사 하강법: 전체 훈련 샘플에 대해 훈련 후 경사 하강 진행

미니 배치 경사 하강법: 전체 훈련 샘플을 작은 훈련 세트인 미니 배치로 나눈 후, 미니배치 훈련 후 경사하강 진행

- 배치 경사 하강법의 경우, 큰 데이터셋에서는 훈련하는데 시간이 많이 들고 결과적으로 경사 하강을 진행하는데 오래걸리기 때문에 작은 훈련세트인 미니 배치로 나뉘어 훈련 후 경사하강을 적용

ex) $m = 7,000,000$ 개의 훈련세트일때 사이즈=1000개인 미니배치 5,000개씩 나눠진법

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} & | & x^{(1001)} & \dots & x^{(2000)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix}$$

$(n_{x,m})$ $x^{1:1000} (n_{x,1000})$ $x^{1001:2000} (n_{x,1000})$ $x^{5000:m} (n_{x,1000})$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} & | & y^{(1001)} & \dots & y^{(2000)} & | & \dots & | & \dots & y^{(m)} \end{bmatrix}$$

$(1,m)$ $y^{1:1000} (1,1000)$ $y^{1001:2000} (1,1000)$ $y^{5000:m} (1,1000)$

* 표기법

- i 번째 훈련 세트: $x^{(i)}$
- l 번째 신경망의 출력: $z^{[l]}$
- t 번째 미니배치: $x^{ft:t}$, $y^{ft:t}$

* Mini-batch gradient descent

for $t=1, \dots, 5000$

① Forward prop on $x^{ft:t}$

$$z^{[1]} = w^{[1]} x^{ft:t} + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$a^{[L]} = g^{[L]}(z^{[L]})$$

} Vectorization
(1000 examples)

② Compute cost $J = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_l \|w^{[l]}\|_F^2$

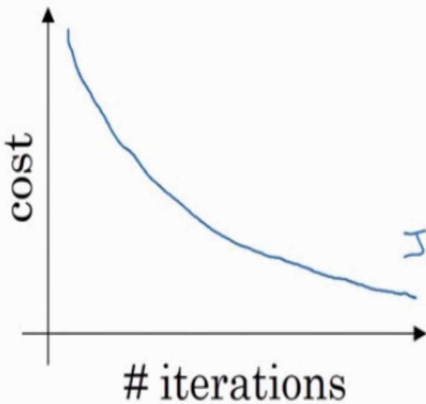
③ Backprop to compute gradient $J^{ft:t}$ (using $(x^{ft:t}, y^{ft:t})$)

$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}, \quad b^{[l]} := b^{[l]} - \alpha db^{[l]}$$

⇒ "1 epoch" 라고 한다.

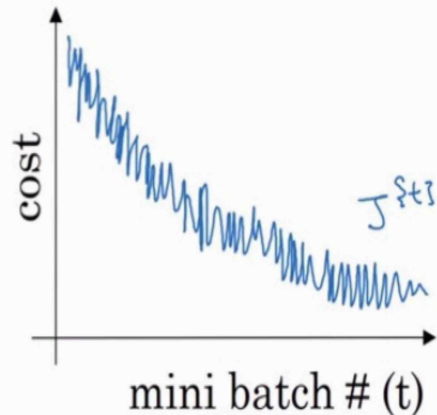
< Understanding Mini-batch Gradient Descent >

Batch gradient descent



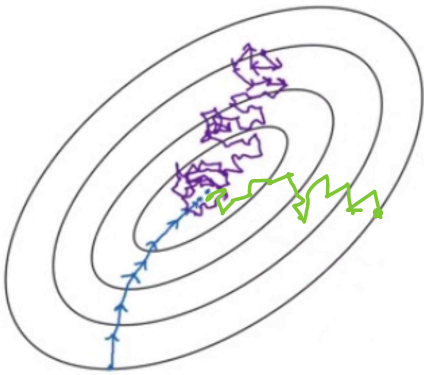
↳ 한번 반복을 돌때마다 비용항수 ↓

Mini-batch gradient descent



↳ 전체적으로 보면 감소하는 형태지만 많은 노이즈를 가짐

* Choosing your mini-batch size



if mini-batch size = m : Batch gradient descent

$$\rightarrow (x^{(i)}, y^{(i)}) = (x, y)$$

⇒ Iteration이 매우 오래 걸림

if mini-batch size = 1 : stochastic gradient descent

$$\rightarrow (x^{(i)}, y^{(i)}) = ((x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)}))$$

⇒ 벡터화하는데에서 독도 노이즈

In practice, somewhere in-between 1 and m

⇒ 가장 빠른 학습 (가장 벡터화, 독도 ↑)

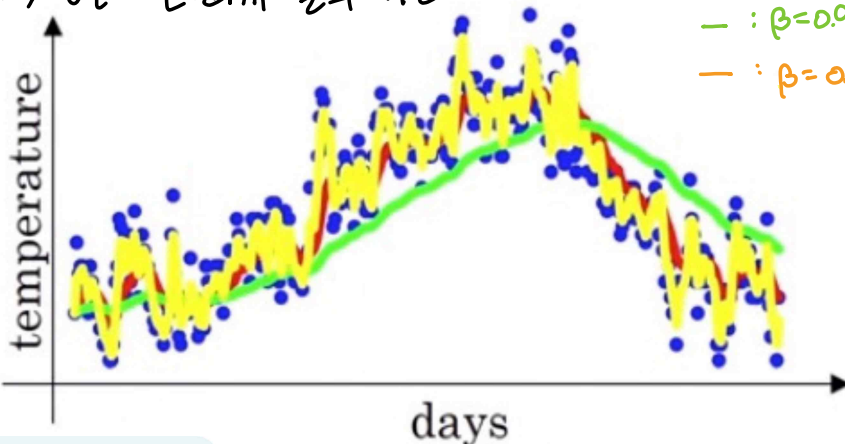
⇒ small training set ($m \leq 2000$): batch gradient descent

typical mini-batch size : 64, 128, 256. 512와 같이 2^n 으로 정함.

< Exponentially Weighted Average >

ex) θ_t : t번째 날의 기온

— : $\beta = 0.9$
— : $\beta = 0.98$
— : $\beta = 0.5$



V_t (exponentially weighted average)

: $\frac{1}{1-\beta}$ 기간 동안의 기온의 평균

$$= \beta V_{t-1} + (1-\beta)\theta_t$$

★ 9월 4일

i) $\beta = 0.9 \approx 10$ day's temperature

ii) $\beta = 0.98 \approx 50$ day's

iii) $\beta = 0.5 \approx 2$ day's

< Understanding Exponentially Weighted Averages >

ex) $V_t = \beta V_{t-1} + (1-\beta) \theta_t$ 의 지수 가중 함수에서, $\beta = 0.9$ 라 하면

$$V_{100} = 0.9 V_{99} + 0.1 \theta_{100}$$

$$V_{99} = 0.9 V_{98} + 0.1 \theta_{99}$$

$$V_{98} = 0.9 V_{97} + 0.1 \theta_{98}$$

...

$$\begin{aligned} \Rightarrow V_{100} &= 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 V_{98}) \\ &= 0.1 \theta_{100} + 0.1 \times 0.9 \theta_{99} + (0.9)^2 V_{98} \\ &= \dots \end{aligned}$$

$$\Rightarrow V_{100} = 0.1 \theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1 \times (0.9)^2 \theta_{98} + \dots$$

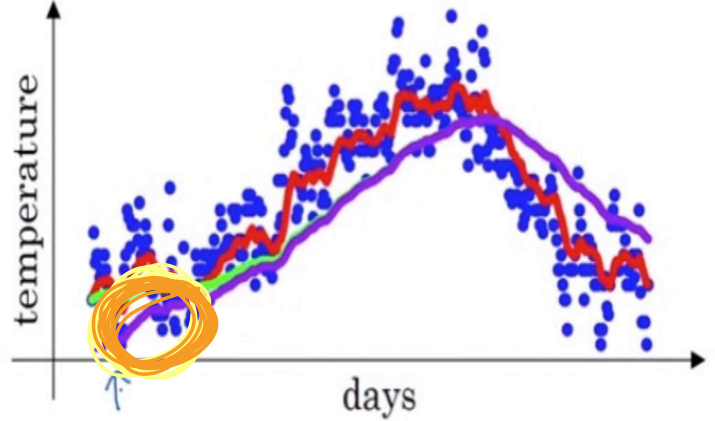
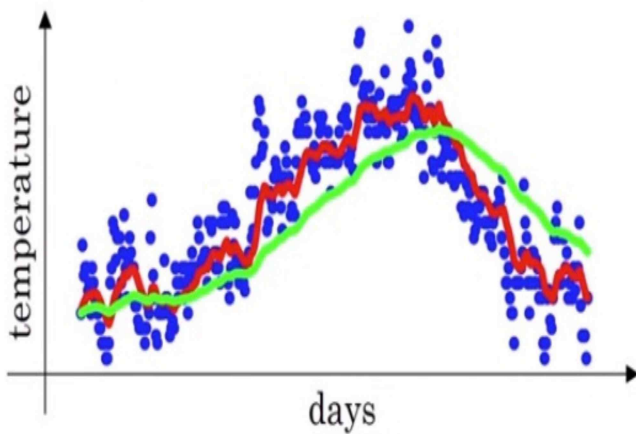
(강조 그래프)

if $\beta = 1 - \epsilon$ 이라고 하면, $(1 - \epsilon)^n = \frac{1}{e}$ 을 만족하는 n 이 n 기간동안 이동하면서 겪는 평균

ex) $\beta = 0.9$ 면 $n = 10$, $\beta = 0.98$ 이면 $n = 50$

- exponentially weighted average의 장점: 아주 적은 메모리를 사용!
($\because V$ 하나의 값만 가지고 있으면 됨)

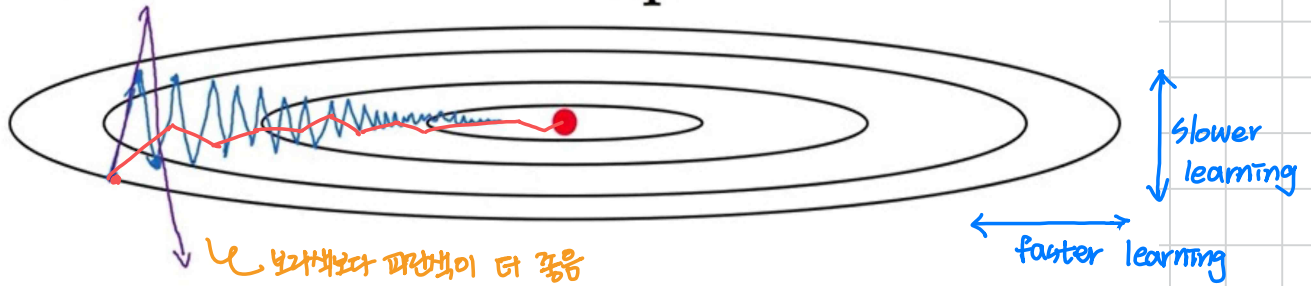
< Bias Correction in Exponentially Weighted Averages >



- 일반적인 지수 가중 평균을 사용하면, 저항비로 표시한 것처럼 초기값에서 차이가 많이 난다.
- 이것을 $\frac{V_t}{1-\beta^t}$ 를 취해서 편향 보정을 해줄 수 있다.
- 그러나 머신러닝에서는 시간이 지나면서 초깃값 그래프가 맞춰지기 때문에 많이 사용하지는 않는다.

< Gradient Descent with Momentum >

Gradient descent example



* Momentum

On iteration t :

Compute dw, db on current mini-batch

$$\rightarrow V_{dw} = \beta V_{dw} + (1 - \beta) dw, \quad V_{db} = \beta V_{db} + (1 - \beta) db \quad \rightarrow \text{논리에서 } (1 - \beta) \text{를 생략한 표기도 자주 사용}$$

$$w := w - \alpha V_{dw}, \quad b := b - \alpha V_{db}$$

Hyperparameter: α, β 이며 $\beta = 0.9$ 를 주로 사용

- 모멘텀은 모멘텀이 없는 경사하강 알고리즘보다 더 잘 작동한다.
- 매 단계의 경사 하강 정도를 더 부드럽게 만들거다 (위의 그림의 빨간색 gradient descent)
- 모멘텀은 이동평균을 적용한다고 생각하면 됨. 즉, 세로축(b)의 경우 위아래 변동폭이 작으면 db 가 작고 가로축(w)의 경우 이동폭이 크기 때문에 dw 가 커진다.

< RMSprop (Root Mean Square prop) >

* RMSprop

On iteration t :

Compute dw, db on current mini-batch

$$\rightarrow S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2 \quad \leftarrow \text{element-wise (요소별 제곱)이며, 자동 평관을 유지}$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \quad \leftarrow \text{일반적인 경사하강을 적용한다면 (세로축 커, 가로축 작, 이동) } dw \text{는 작고 } db \text{는 큰 값 많거림}$$

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}, \quad b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}} \quad \leftarrow \epsilon = 10^{-8}$$

- RMSprop는 이동폭이 큰 곳 (db)에서는 큰 값 (S_{db})로 나누어주고, 이동폭이 작은 곳 (dw)에서는 작은 값 (S_{dw})로 나누어주어서, 각각 기온의 하중효과와 작은값, 큰값으로 업데이트되기 때문에 전동을 줄이는데에 도움을 주고, 더 빠르게 수렴하는 효과를 가짐

<Adam Optimization Algorithm>

* Adam (= Adapted Momentum estimation)

- 광범위한 딥러닝 알고리즘에서 잘 작동하여 'momentum + RMSprop' 방법

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

On iteration t :

Compute dw, db using current mini-batch

$$\rightarrow V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \quad \leftarrow \text{"momentum"}$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \quad \leftarrow \text{"RMS prop"}$$

$$V_{dw}^{correct} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{correct} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{correct} = S_{dw} / (1 - \beta_2^t), \quad S_{db}^{correct} = S_{db} / (1 - \beta_2^t)$$

$$W := W - \alpha \frac{V_{dw}^{correct}}{\sqrt{S_{dw}^{correct} + \epsilon}}, \quad b := b - \alpha \frac{V_{db}^{correct}}{\sqrt{S_{db}^{correct} + \epsilon}}$$

- Hyper parameters choice

- α : needs to be tune

- β_1 : 0.9 (=dw)

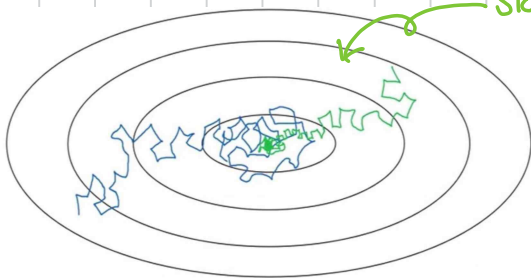
- β_2 : 0.999 (=dw²)

- ϵ : 10^{-8}

<Learning Rate Decay>

- 학습률 감소: 시간이 지날수록 점점 학습률을 작게 해서 최적값을 더 빨리 찾도록 함

slowly reduced α



⇒ 작은 미니배치일수록 noise가 심해지기 때문에, 학습률(α)가 동일하면 최적값에 부정확히 힘들다.

∴ α 를 점점 줄여가면서 최적값을 찾도록

• 1 epoch: 1 pass through all data

$$\alpha = \frac{1}{1 + \text{decay rate} \times \text{epoch num}} \alpha_0$$

ex $\alpha_0 = 0.2$, decay rate = 1

epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
...	...

→ 점진적으로 감소

* Other learning rate decay methods

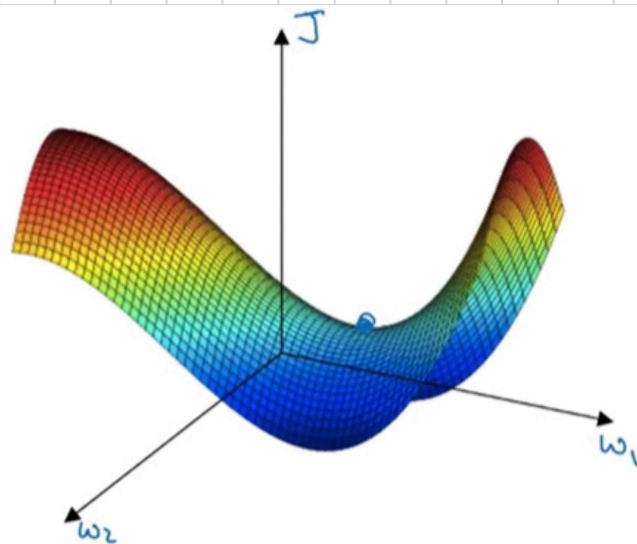
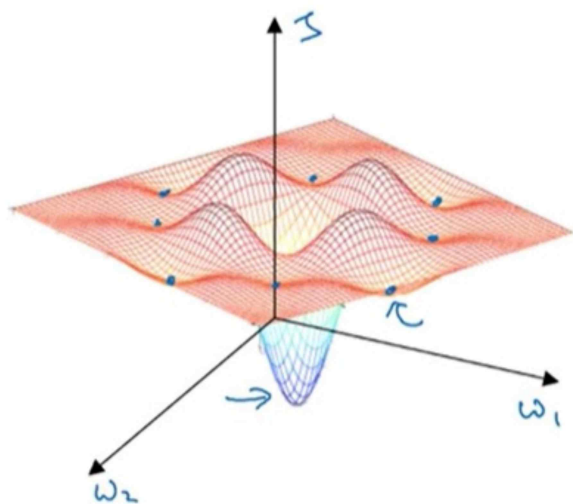
① $\alpha = 0.95^{\text{epoch num}}$ α_0 : exponential decay

② $\alpha = \frac{k}{\sqrt{\text{epoch num}}} \alpha_0$

③ $\alpha = \frac{k}{\sqrt{\text{batch num}}} \alpha_0$

< The Problem of Local Optima >

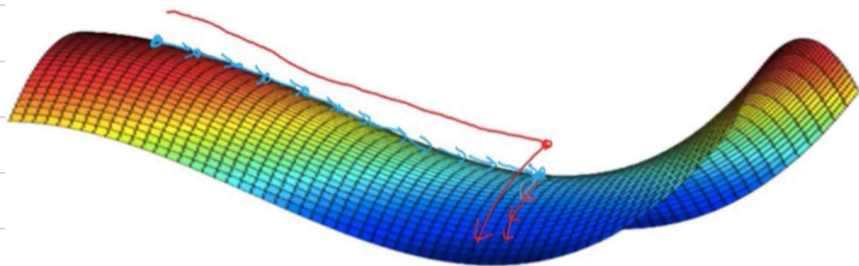
*



왼쪽 사진에서처럼 local optima problem이 생길 수 있지만,

산정양에서 매우 큰 hyper parameter를 다루고 있으므로 이러한 문제에 빠질일은 거의 x 대부분의 경우에 오른쪽처럼 global optima로 회귀한다.

* Problem of plateaus



plateau : 경사가 완만해지면서 gradient descent를 하는 속도가 굉장히 느려지는 것