

Optimization Algorithms

Mini batch Gradient Descent

- $n = 5,000,000$ 개의 dataset이 있다고 가정 하면.

이것을 한 번에 처리하기에는 엄청난 처리 시간이 필요.

- 시간이 걸리면 mini-batch 활용하면 시간과 같이省

$$X^{(1)} = [x^{(1)}, x^{(2)}, \dots, x^{(m)}] \rightarrow \text{节省},$$

$$\left. \begin{array}{l} X^{(2)} = [x^{(1w)}, \dots, x^{(2w)}] \\ \vdots \end{array} \right\}$$

$$X^{(5000)} = [x^{(1,000,001)}, \dots]$$

- Y 도 이와 같이 1000개씩 잘라준다.

- mini-batch $\Rightarrow t$; $X^{(t)}, Y^{(t)}$

for $t = 1, \dots, 500$

AL, caches = forward-prop ($X^{(t)}, Y^{(t)}$)

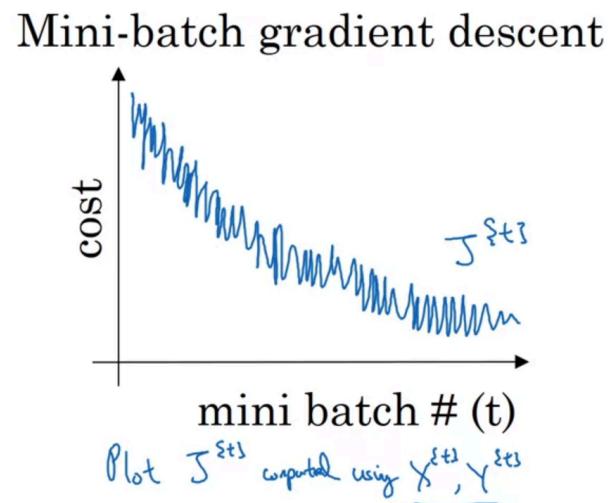
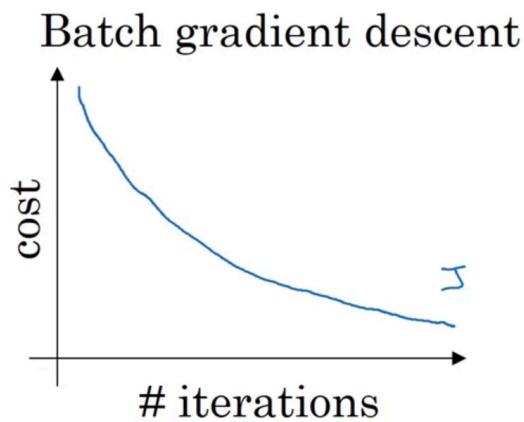
cost = compute-cost (AL, $Y^{(t)}$)

grads = backward-prop (AL, caches)

update-parameters (grads).

- mini-batch을 사용하는 것이 속도 면에서 훨씬 빠르다.
작동된다.

Understanding mini batch gradient descent



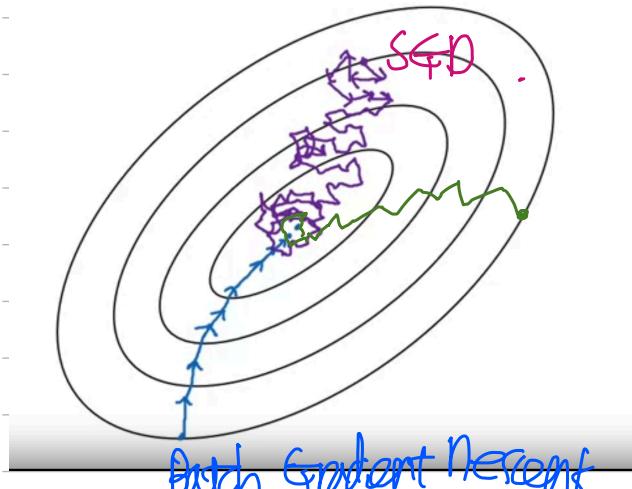
- Gradient Descent에서는 각 반복마다 J 가 감소.
- Mini-batch gradient descent에서는 각 다른 dataset을 사용하기 때문에 J 가 감소하는 경향을 small noise라고 한다.

• Mini-batch Size

mini-batch size = $m \Rightarrow$ Batch Gradient Descent

mini-batch size = 1 \Rightarrow stochastic Gradient Descent

mini-batch size = between 1 and $m \Rightarrow$ Mini-batch Gra...



Batch Gradient Descent.

\rightarrow Too long per iteration

SGD

\rightarrow TOO noisy regarding cost minimization.

- Won't ever converge (reach the minimum cost)
- lose speedup from vectorization

Mini-batch Gradient Descent

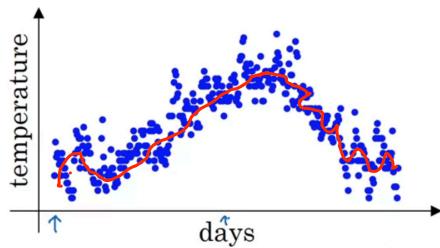
\rightarrow Fastest learning \Rightarrow vectorization advantage. make process without waiting to process the entire training set

* training set의 모든 데이터는 Batch Gradient Descent, 다른 mini-batch 크기 ($2^n, 64, 128, \dots$)

Exponentially weighted averages

Temperature in London

$$\begin{aligned}\theta_1 &= 40^{\circ}\text{F} \quad 4^{\circ}\text{C} \\ \theta_2 &= 49^{\circ}\text{F} \quad 9^{\circ}\text{C} \\ \theta_3 &= 45^{\circ}\text{F} \quad \vdots \\ &\vdots \\ \theta_{180} &= 60^{\circ}\text{F} \quad 15^{\circ}\text{C} \\ \theta_{181} &= 56^{\circ}\text{F} \quad \vdots \\ &\vdots\end{aligned}$$



$$V_0 = 0$$

$$V_1 = 0.9V_0 + 0.1\theta_1$$

$$V_2 = 0.9V_1 + 0.1\theta_2$$

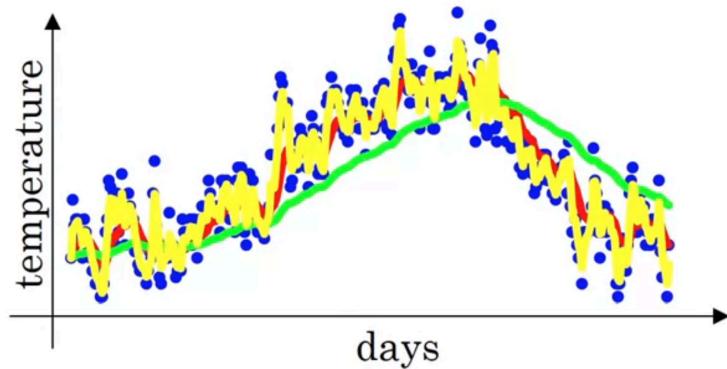
\vdots

temperature의 값

Noise가 있다.

noise 없애기 위해 average 사용. $V_t = 0.9V_{t-1} + 0.1\theta_t$.

식에 포함된 0.9는 smoothness의 정도 결정



V_t can be considered to be approximate average over $\frac{1}{1-\beta}$ days

temperature

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

$\beta = 0.9 \rightarrow 10 \text{ days}$: 현재 sample을 포함, 과거 9개의 sample을 고려한다는 뜻

$$\beta = 0.98 \rightarrow 50 \text{ days}$$

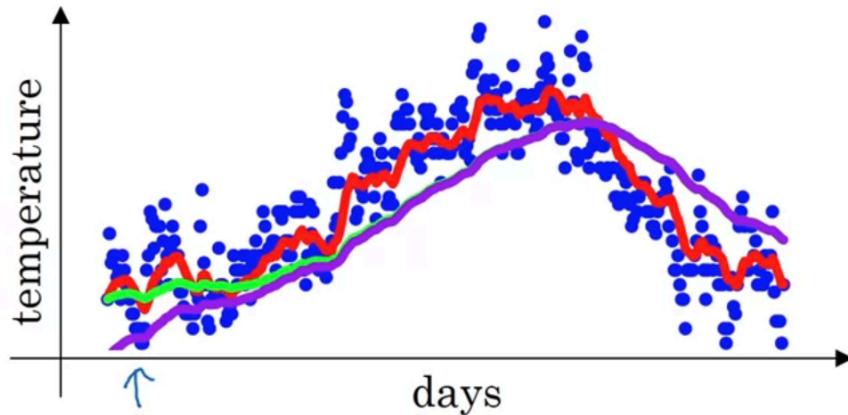
$$\beta = 0.5 \rightarrow 2 \text{ days}$$

β 를 너무 카워 너무 smooth해지면 기준의 peak가 뒤로 밀리는 delay 발생 가능.
이런 delay를 감내할 수 있는 분야에서 값을 전한다.

$$V_{100} = 0.1\theta_{100} + 0.1\theta_{99} + 0.1(0.9)^2\theta_{98} + 0.1(0.9)^3\theta_{97} + \dots$$

$\Rightarrow \beta$ 이 exponential하게 감소하기에 exponentially weighted 된다.

Bias correction in exponentially weighted averages



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

$$V_0 = 0$$

$$V_1 = 0.98 V_0 + 0.02 \theta_1 = 0.02 \theta_1$$

$$V_2 = 0.98 V_1 + 0.02 \theta_2 = 0.98 \times 0.02 \theta_1 + 0.02 \theta_2$$

$$= 0.0196 \theta_1 + 0.02 \theta_2$$

$$V_t \rightarrow \frac{V_t}{1-\beta^t} \quad t=2 \Rightarrow 1-\beta^2 = 1-0.98^2 = 0.0396$$

$$\frac{V_2}{0.0396} = 0.49 \theta_1 + 0.51 \theta_2$$

• V_1 을 구할 때 V_0 이 0인가 때문이 일으킨 정의착오.

하지만 초기값 V_0 은 0이기 때문에 V_1 을 구할 때 $V_0=0$ 을 사용해버리면 V_1 의 값이 0으로 가까워지는 문제점 발생.

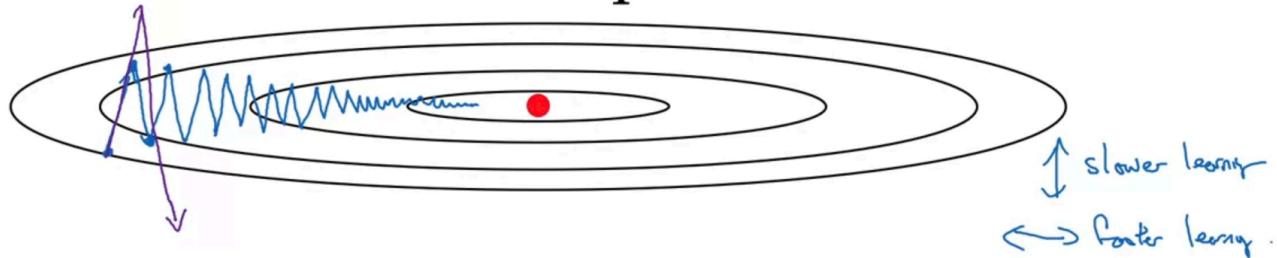
• Bias Correction $V_t \rightarrow \frac{V_t}{1-\beta^t}$ 는 starting point에서.

초기선과 끝선을 맞추시하기 위한 방법.

t 가 엄청 커지면, β^t 가 0에 수렴하기 때문에 V_t 가 0이 가까워지고 자연스럽게 bias를 고려하지 않기 된다.

Gradient descent with momentum

Gradient descent example



Momentum : On iteration t :

Compute dW and db on current mini-batch.

$$V_{dW} = \beta V_{dw} + (1-\beta) dW$$

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

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

Hyperparameters : α, β

$$\beta = 0.9$$

average over last ≈ 10 gradients

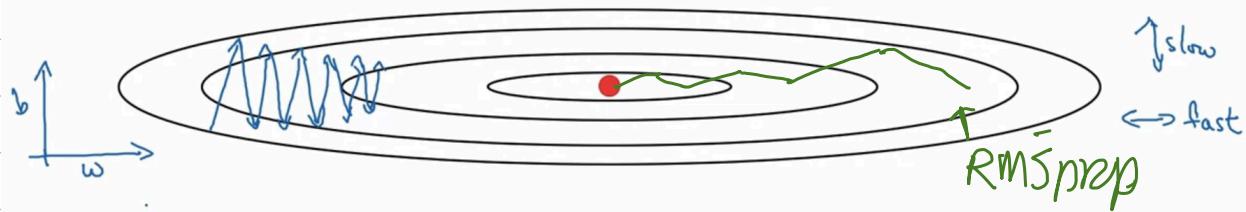
mini-batch 을 통해 gradient descent 를 사용하면 . 빨간 점이

도달하는 데 굉장히 무서워지게 느껴질 것이다.

그러나 noise 가 덜 차오록 $(1-\beta)$ 을 smooth 를 해준다.

RMSprop - Root Mean Square prop.

RMSprop



- RMSprop은 단순히 smooth한 경로가 아니라 gradient minimum을 가는 속도를 조절한다. (세로방향의 gradient는 slow하면 좋겠고, 가로방향의 gradient는 fast하면 좋겠다고 생각).
- On Iteration t :

Compute dW and db on current minibatch.

$$S_{dW} = \beta S_{dW} + (1-\beta) dW^2$$

want small S_{dW}

$$S_{db} = \beta S_{db} + (1-\beta) db^2$$

want large S_{db} .

- RMSprop은 W 를 커하고, b 를 줄이고 싶다.

S_{dW} 가 1보다 작으면 W 의 값이 커지고, S_{db} 가 1보다 크면.

b 의 값이 커지게 된다. 작은 S_{dW} 와 S_{db} 가 필요.

$$W := W - \alpha \frac{dW}{\sqrt{S_{dW}}}$$

$$b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

Adam optimization

- Adam은 $V_{dw} = 0$ or $S_{dw} = 0$, $V_{db} = 0$, $S_{db} = 0$

On iteration t :

Compute dW, db on current minibatch.

$$\text{“} V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dW, V_{db} = \beta_1 V_{db} + (1-\beta_1) db \text{”} \Rightarrow \text{momentum with } \beta_1$$

$$\text{“} S_{dw} = \beta_2 V_{dw}^2 + (1-\beta_2) dW^2, S_{db} = \beta_2 V_{db}^2 + (1-\beta_2) db^2 \text{”} \Rightarrow \text{RMSprop with } \beta_2$$

$$\text{“} \overset{\text{corrected}}{V_{dw}} = \frac{V_{dw}}{(1-\beta_1^t)}, \overset{\text{corrected}}{V_{db}} = \frac{V_{db}}{(1-\beta_1^t)} \text{”}$$

$$\text{“} \overset{\text{corrected}}{S_{dw}} = \frac{S_{dw}}{(1-\beta_2^t)}, \overset{\text{corrected}}{S_{db}} = \frac{S_{db}}{(1-\beta_2^t)} \text{”}$$

$$W := W - \alpha \frac{\overset{\text{corrected}}{V_{dw}}}{\sqrt{\overset{\text{corrected}}{S_{dw}}} + \epsilon}, b := b - \alpha \frac{\overset{\text{corrected}}{V_{db}}}{\sqrt{\overset{\text{corrected}}{S_{db}}} + \epsilon}$$

- Adam은 v 는 그때 gradient 사용, s 는 dW^2 사용.

W 와 b 를 구할 때는 gradient가 빠른 값을 유지하도록.

분포로 나누어진다.

Hyper parameters α : needs to be tuned.

recommended by default. $\left\{ \begin{array}{l} \beta_1 = 0.9 \text{ (weighted average for } dW) \\ \beta_2 = 0.99 \text{ (weighted average for } dW^2) \\ \epsilon = 10^{-8} \end{array} \right.$

Learning rate decay

1 epoch = 1 pass through data.

$$\alpha = \frac{\alpha_0}{1 + \text{decay_rate} \times \text{epoch_num}}$$

$\alpha = 0.2$,
 $\text{decay_rate} = 1$

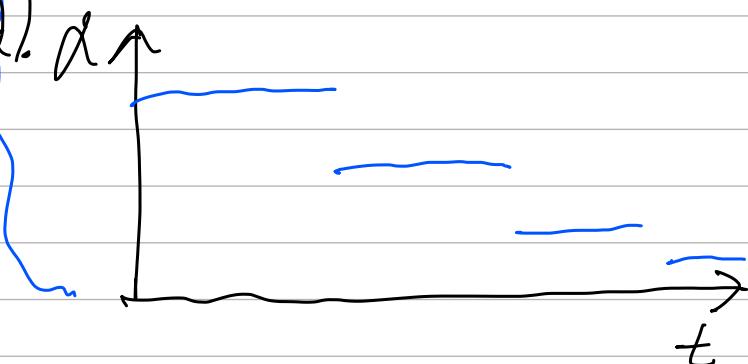
Epoch	α
1	0.1
2	0.07
3	0.05
:	:
:	:

α 와 decay_rate는 hyper parameter.

iteration의 증가수로 learning rate α 는 감소.

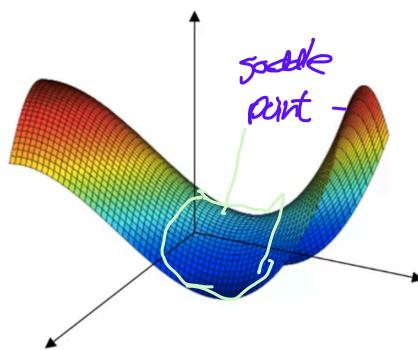
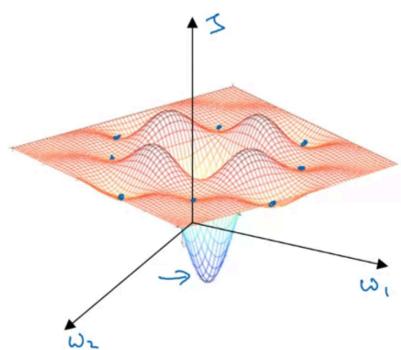
$\alpha = 0.95^{\text{epoch_num}} \cdot \alpha_0$: exponentially decaying

Formula
(based.) $\alpha = \frac{k}{\sqrt{\text{epoch_num}}} \cdot \alpha_0$ or $\alpha = \frac{k}{\sqrt{t}} \cdot \alpha_0$



α 값을 epoch에 따라 일정하게 조절하기도 한다.

The problem of local optima



과거에는 왼쪽 그림처럼 solution space가 엄청 복잡할 것이다 - 생각함.

하지만 실제로는 오른쪽 그림처럼 딸린 장소들 평평하게 생긴

solution space 존재.

- plateaus can make learning slow
 - Plateau is a region where the derivative is close to zero for a long time
 - This is where algorithms like momentum, RMSprop or Adam can help -

$$0,5V_0 + 0,5\theta_1$$

$$V_1 = 15 \cdot 0,5 V_1 + (1-\beta) V_2 \\ 15 = 15 - 15$$

$$\underline{V_2 = 15} //$$