

유런 10주차 정리

미니 배치 경사 하강법

모델을 빠르게 학습시키는 방법

큰 데이터 세트에서 학습시키는 것은 오래 걸리는 것이기에 좋은 최적화 알고리즘 방법을 찾아야한다.

벡터화 하면 명시적인 for 문이 없어도 훈련이 가능하다.

m 이 너무 크다면 똑같이 느낄 수 있다. 경사 하강법의 다음 단계를 밟기 전에 모든 training set 를 밟아야한다. 거대한 훈련 샘플을 다 돌기 전에 경사 하강법이 진행되도록 하면 더 빠른 최적화를 이룰 수 있다.

Batch vs. mini-batch gradient descent

Vectorization allows you to efficiently compute on m examples.

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

(n_x, m) $X^{\{1\}} (n_x, 1000)$ $X^{\{2\}} (n_x, 1000)$ $X^{\{5,000\}} (n_x, 1000)$

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

$(1, m)$ $Y^{\{1\}} (1, 1000)$ $Y^{\{2\}} (1, 1000)$ $Y^{\{5,000\}} (1, 1000)$

What if $m = 5,000,000$?

5,000 mini-batches of 1,000 each

Mini-batch t : $x^{(i)}$
 $\begin{bmatrix} x \\ z \\ y \end{bmatrix}$ $\{t\}$

다음 슬라이드에 나와있는데
전체 훈련 세트 X, Y를 한 번에 진행시키지 않고

Andrew Ng

baby training set(5000000개) = mini batch(1000개) $x^{\{1000\}}$

$x^{\{t\}}, y^{\{t\}}$ 를 가지고 경사하강법을 진행한다.

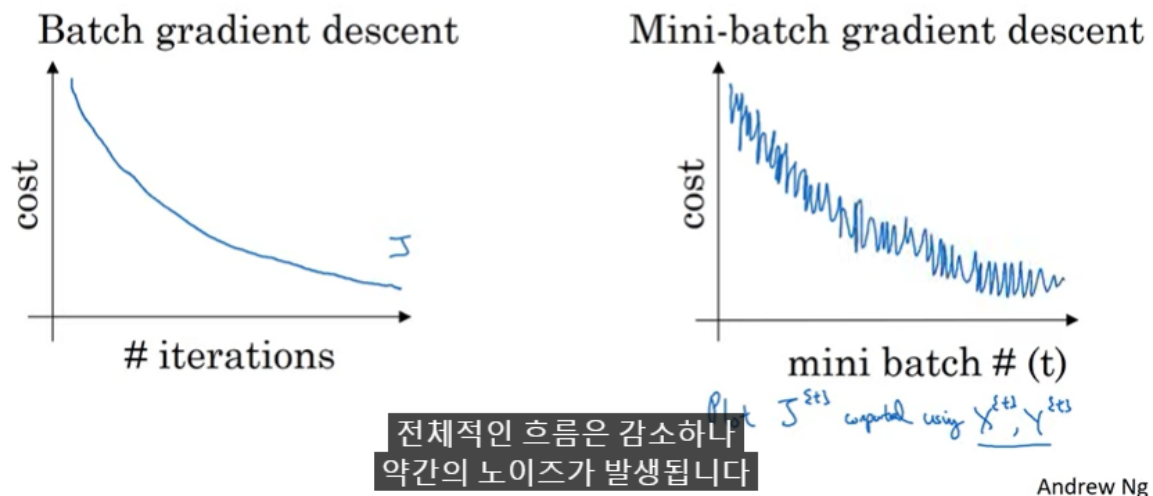
따라서 배치 경사 하강법에서 훈련 세트를 거치는 한 반복은 오직 하나의 경사 하강법을 거치고,

미니 배치 경사 하강법의 경우, 훈련 세트를 거치는 한 반복은 5000개의 경사 하강 단계를 거친다.

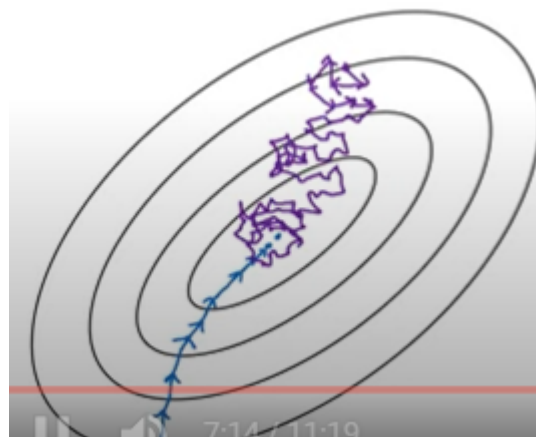
훈련 세트가 많다면 미니 배치 경사 하강법이 빠르게 작용한다.

미니 배치 경사하강법 이해하기

Training with mini batch gradient descent



모든 반복에서 아래로 내려가진 않지만 큰 흐름은 내려간다.



파란색이 batch, 보라색이 mini-batch

보통 미니 배치 크기는 2의 제곱수로 선택한다.

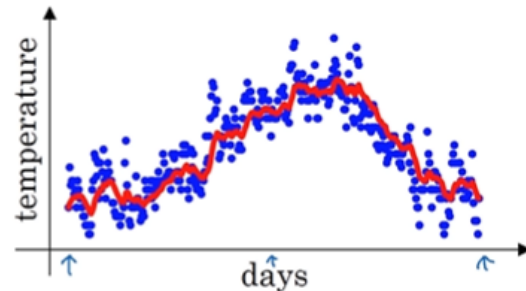
모든 배치에 대해 cpu, gpu 메모리에 맞는지 확인해본다. 메모리에 맞지 않으면 성능이 갑자기 나빠지고 훨씬 나빠지게 된다.

따라서 미니 배치 사이즈를 어떻게 선택하는지에 따라 학습 속도의 차이가 나기에 최적의 값을 찾아내는 것이 중요하다.

지수 가중 이동 평균

Temperature in London

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



$$\begin{aligned}v_0 &= 0 \\ v_1 &= 0.9 v_0 + 0.1 \theta_1 \\ v_2 &= 0.9 v_1 + 0.1 \theta_2 \\ v_3 &= 0.9 v_2 + 0.1 \theta_3 \\ &\vdots \\ v_t &= 0.9 v_{t-1} + 0.1 \theta_t\end{aligned}$$

익번 기온이 지수가중평균을 얻게 됩니다

이전 값에 0.9를 곱하고 해당 날의 기온에 0.1를 곱한다. 과거의 경향성과 새로운 경향성을 더한다.

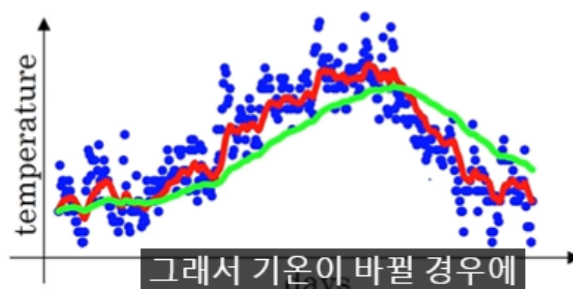
Exponentially weighted averages

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

$\beta = 0.9$: ≈ 10 days' temperat.
 $\beta = 0.98$: ≈ 50 days

v_t is approximately
average over
 $\approx \frac{1}{1-\beta}$ days'
temperature.

$$\frac{1}{1-0.98} = 50$$



그래서 기온이 바뀔 경우에
지수가중평균 공식은 더 느리게 적응합니다

Andrew P

v_t 는 $1/(1-\beta)$ 기간 동안 기온의 평균을 의미

β 가 0.9 와 같은 경우 이는 10일 동안 기온의 평균과 같다. (붉은색)

β 가 1 과 가까우면 50일의 평균과 같다. 값이 커질수록 곡선이 부드러워진다(연두). 더 많은 날짜의 평균을 사용하기 때문이다.

β 가 0.5 이면 2일의 평균을 쓰기 때문에 더 크게 진동한다.

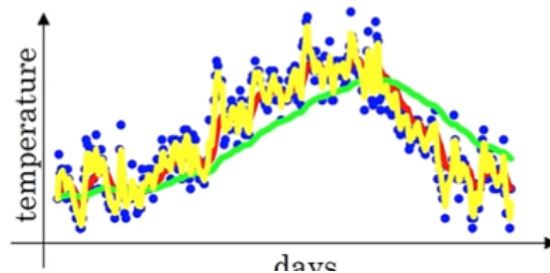
Exponentially weighted averages ^{moving}

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

$\beta = 0.9$: ≈ 10 days' temper
 $\beta = 0.98$: ≈ 50 days
 $\beta = 0.5$: ≈ 2 days

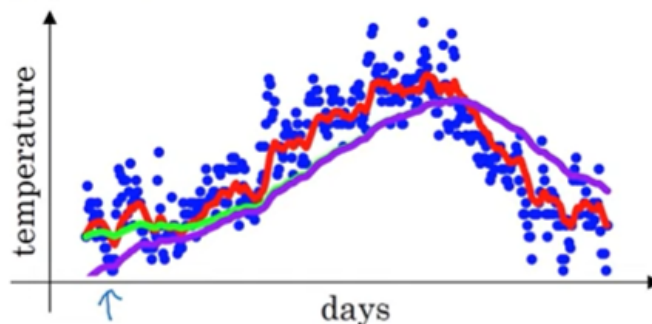
V_t is approximately
 average over
 $\rightarrow \approx \frac{1}{1-\beta}$ days' temperature.

$$\frac{1}{1-0.98} = 50$$



- 얼마의 기간이 이동하면서 평균이 구해졌는가? 는 아래의 식으로 대략적으로 구할 수 있다.
 - $\beta = (1-\epsilon)$ 라고 정의 하면
 - $(1-\epsilon)^n = 1/e$ 를 만족하는 n 이 그 기간이 되는데, 보통 $\epsilon=1$ 으로 구할 수 있다.
- 지수 가중 이동 평균의 장점은 구현시 아주 적은 메모리를 사용한다는 것이다.

Bias correction



$$\beta = 0.98$$

보라색 bias 를 바로잡을 것이다.

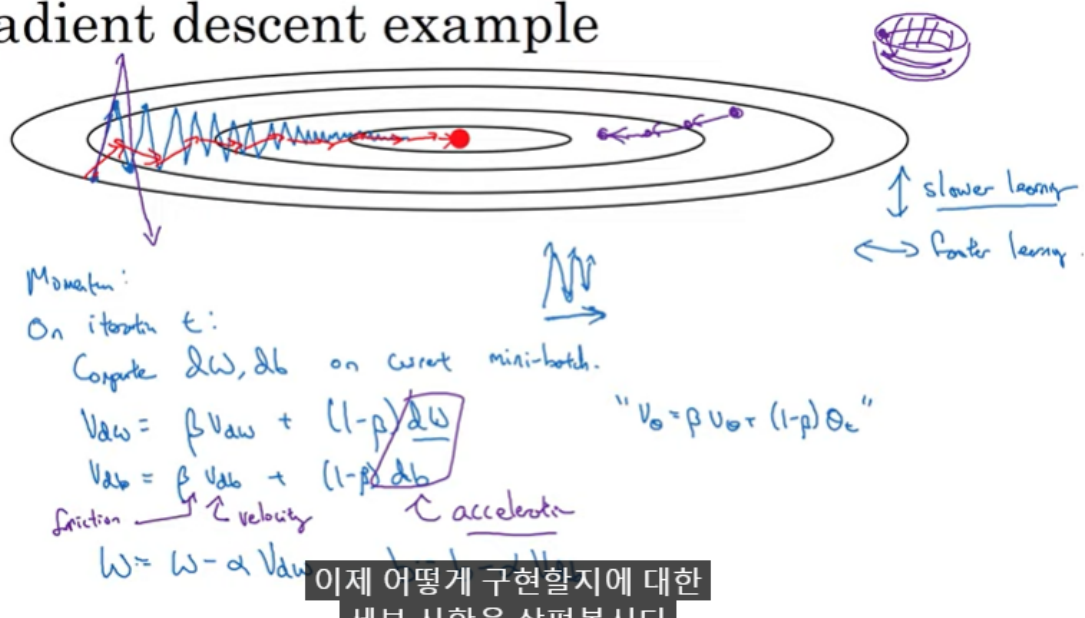
보라색 선의 초기 단계에서 0에 가까운 것을 알 수 있다. 한 해의 첫 두 날짜를 추정한 값의 좋지 못한 결과이다.

- 따라서 $v_t / (1 - \beta_t)$ 를 취해서 초기 값에서 실제값과 비슷해지게 한다.
- 보통 머신러닝에서 구현하지는 않는다. 시간이 지남에 $(1 - \beta_t)$ 는 1에 가까워져서 원하는 값과 일치하게 되기 때문이다.

Momentum 최적화 알고리즘

이전 값을 보고 결정한다. 수직 방향은 양수 음수가 번갈아 가면서 나오므로 잘 안 움직이고 수평 방향은 오른쪽을 향하기 때문에 오른쪽으로 가속화된다.

Gradient descent example



Algorithm

Implementation details

$$v_{dw} = 0, v_{db} = 0$$

On iteration t :

Compute dW, db on the current mini-batch

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

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

$$W = W - \alpha v_{dw}, b = b - \alpha v_{db}$$



Hyperparameters: α, β

모멘텀이 있는 경사 하강법에 대한
논문을 읽어보면

loss % 10 grads

Andre

RMSProp 최적화 알고리즘(Root Mean Square)

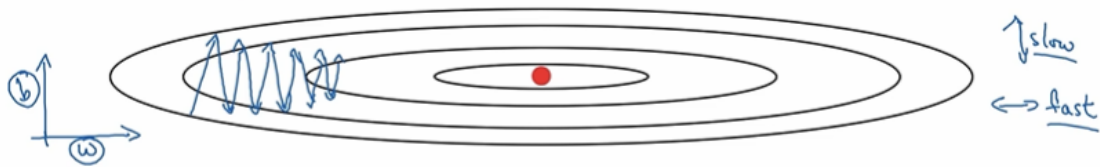
반복 t 에서 현재의 미니배치에 대한 보통의 도함수 dw, db 를 계산할 것이다.

지수가중 평균을 유지하기 위해 새로운 표기법 s_{dw} 을 사용할 것이다

⇒ 도함수의 제곱을 지수가중평균하는 것이다. 비슷하게 s_{db} 도 사용할 것이다.

매개변수 업데이트 식은 아래와 같다.

RMSprop



On iteration t :

Compute dw, db on current mini-batch

$$S_{dw} = \beta S_{dw} + (1-\beta) \frac{dw^2}{\text{element-wise}}$$

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

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw}}}$$

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

따라서 이 두 항, s_{dw} 와 s_{db} 에서
우리가 원하는 것은

Adam 최적화 알고리즘

Adam optimization algorithm

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

On iteration t :

Compute dw, db using current mini-batch

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

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

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

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

$$w := w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}} + \epsilon}}$$

$$b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$

따라서 이 알고리즘은 많은
하이퍼파라미터가 있습니다

Hyperparameters choice:

→ α : needs to be tune
→ β_1 : 0.9 → (\underline{dw})
→ β_2 : 0.999 → ($\underline{dw^2}$)
→ ϵ : 10^{-8}

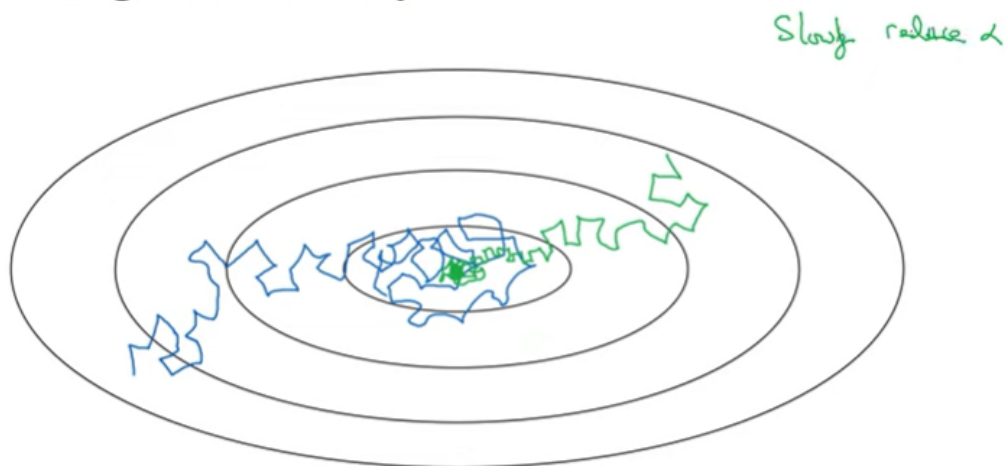
Adam: Adaptive moment estimation

학습률 감쇠

왜 학습률 감쇠가 필요한가

- 작은 미니 배치 일수록 minimum 에서 배회할 가능성이 높기 때문이다.

Learning rate decay



학습률 감쇠 방법

- $\alpha = 1 + (\text{decay rate} \times \text{epoch num}) \alpha_0$
- $\alpha = 0.95^{\text{epoch num}} \alpha_0$ (exponential decay 라고 부른다.)

- $k/\sqrt{epochnum})\alpha_0$
- $k/\sqrt{batchnum})\alpha_0$
- step 별로 α 다르게 설정