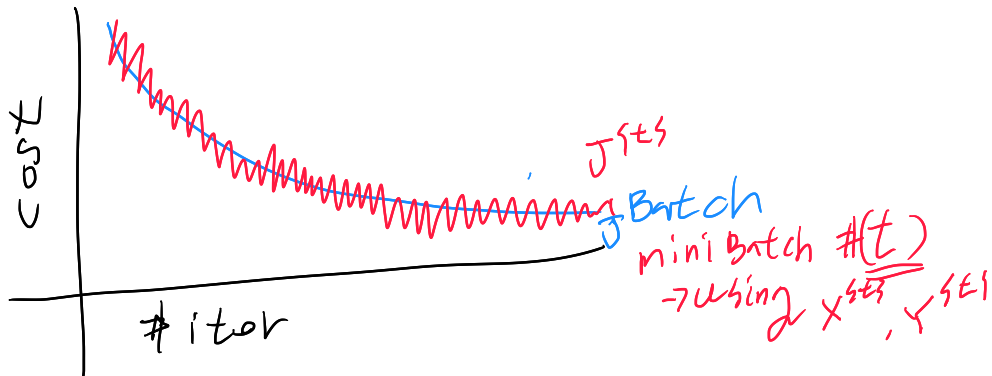


batch \rightarrow SGD 에 전체 데이터 삽입

mini-batch \rightarrow SGD 에 일부 데이터만

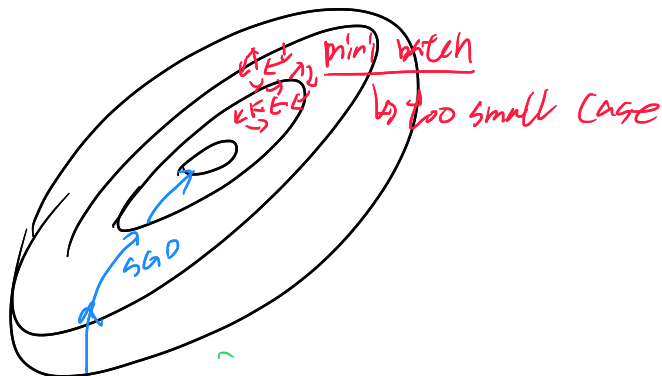
- Y 를 가지고 학습한 것과 동일할 것.

How it work?



Choosing mini-batch size.

- Size down 은 local minimum 으로 갈 확률을 높힘 Loss variation
- Huge Size 는 too Slow 의 소지가 있음 too Long per iter.
- Not too big or too small



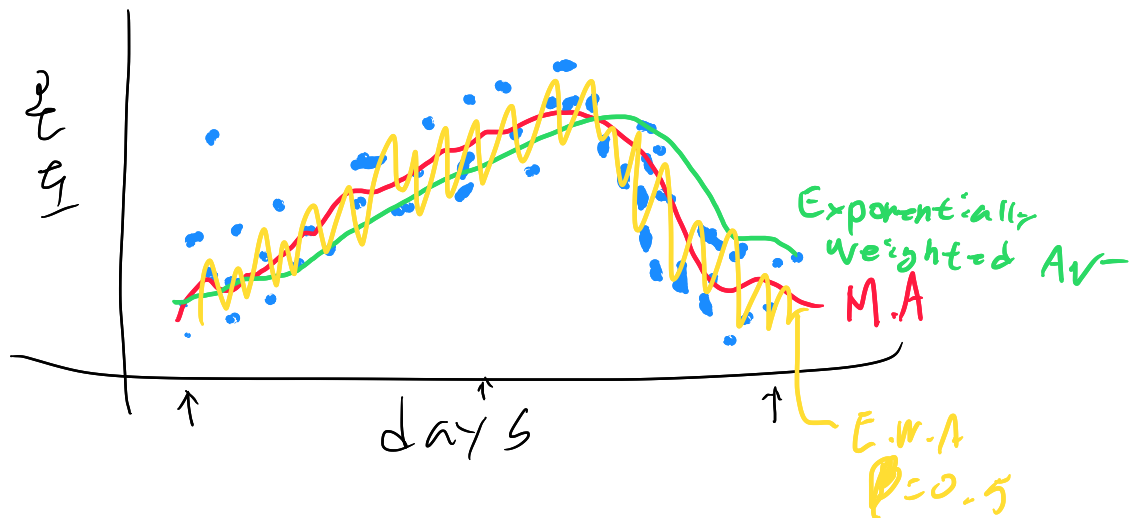
If small training set : use batch

($m < 2000$)

Typical mini-batch Size:

64, 128, 256, 512

Cpu. Gpu 메모리 허용 범위에서 크게.



$$\begin{aligned} V_0 &= 0 \\ V_1 &= 0.9V_0 + 0.1\theta_1 \\ V_2 &= 0.9V_1 + 0.1\theta_2 \end{aligned} \quad \left. \vphantom{\begin{aligned} V_0 &= 0 \\ V_1 &= 0.9V_0 + 0.1\theta_1 \\ V_2 &= 0.9V_1 + 0.1\theta_2 \end{aligned}} \right\} \text{moving Average}$$

β ← hyper parameter.

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

$\beta = 0.9$: ≈ 10 days tm.

$\beta = 0.99$: ≈ 50 days tm.

$\beta = 0.5$: ≈ 2 days tm.

more slowly
↓

Ex -
we -
Av -

Exponentially weighted averages

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

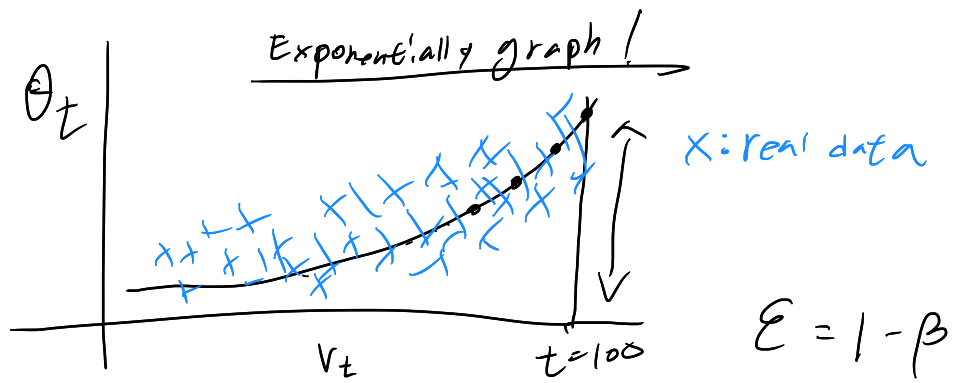
$$v_{100} = 0.9v_{99} + 0.1\theta_{100}$$

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

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

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

So!



구현 이슈

$V = 0$

- $V = \text{Beta} \times V + (1 - \text{Beta}) \times \text{Theta}$

메모리 절약 가능함.

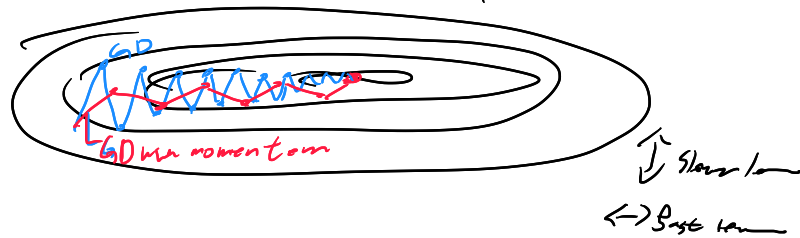
Bias correction

$$\frac{V_t}{1 - \beta^t}$$

$$t = 2; 1 - \beta^t = 1 - (0 - 98)^2 = 0.396$$

Gradient descent with momentum

example of GR - de



momentum:

on iteration t :

compute dW, db on mini-batch

$$\begin{aligned} v_{dW} &= \beta v_{dW} + (1-\beta) dW \\ v_{db} &= \beta v_{db} + (1-\beta) db \end{aligned} \quad \leftarrow \text{accelerator}$$

Implementation details

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

On iteration t :

Compute dW, db on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1-\beta) dW \quad \left| \quad v_{dW} = \beta v_{dW} + dW \leftarrow$$

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

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

$$\frac{v_{dW}}{1-\beta^t}$$

Hyperparameters: α, β

$$\beta = 0.9$$

average over loss & 10 gradients

Andrew Ng

RMSprop

MLP

GD

GD w...

↓ Slower

↔ Faster

Compute $d_{w, db}$ on mini-batch & element-wise

$$S_d b = \beta S_d b + (1-\beta) d b^2 \ll \text{large}$$

$$b \pm = b - d \frac{d' b}{\sqrt{5} b}$$

$$Vdw = 0, Sdw = 0$$

Sdw is Rms prop

$$S_{dW}^{corrected1} = S_{dW} / (l - \beta_2^t)$$

 $\beta_1 : 0.9 (dw)$

$$\beta_2 : 0 - 999 \text{ (} dW^2 \text{)}$$

Learning rate decay

1 epoch = 1 pass through data

1 epoch = 1 pass through data.

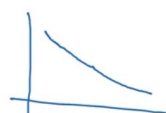
$$\alpha' = \frac{1}{1 + \text{decay-rate} * \text{epoch-num}} \alpha_0$$

E_{podn}	α
1	0.1
2	0.67
3	0.5
4	0.4



$$d_0 = 0.2$$

$$\text{decay rate} = 1$$



rate 가 너무 크면, 최적해를 못하고
너무 작으면, 너무 오래 걸림.


Other learning rate decay methods

formula

$$\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \quad - \text{exponentially decay.}$$

$$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \cdot \alpha_0$$

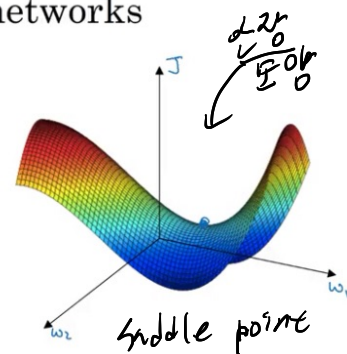
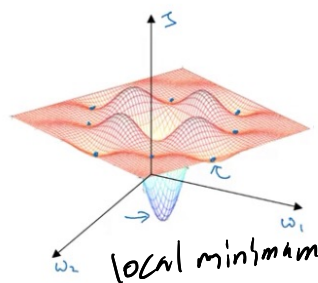
discrete staircase



Manual decay.

Local optima

Local optima in neural networks



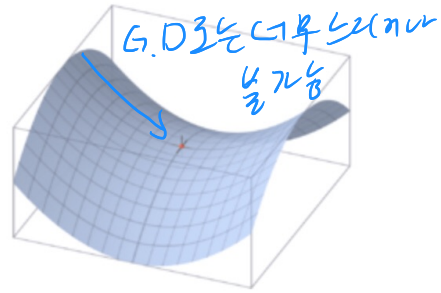
안장점

☆



안장점(鞍裝點; saddle point)은 **다변수 실함수**의 변역에서, 어느 방향에서 보면 **극대값**이지만 다른 방향에서 보면 **극소값**이 되는 점이다.

수학에서 안장점이란 **정류점**이지만 **극값**을 가지지 않는 점을 말한다. 이차원의 시각에서 어느 방향에서 보면 아래로 굽어있지만 다른 방향에서 보면 위로 굽어있는 전형적인 모양이 고개 모양 혹은 안장과 닮았다고 하여 붙여졌다. **등고선**의 관점에서 이차원에서의 안장점은 두 등고선의 교차점으로 나타난다.



함수 $f(x, y) = x^2 - y^2$ 의 안장점.

≡ 목차 ▾

^ 정의



점 (a_1, \dots, a_n) 이 다변수 실함수 $f(x_1, \dots, x_n)$ 의 안장점이라는 것은, **영벡터**가 아닌 2개의 **벡터** (M_1, \dots, M_n) 와 (m_1, \dots, m_n) 가 **존재**하고 그 두 벡터에 대하여,

함수 $g(t) = f(a_1 + tM_1, \dots, a_n + tM_n)$ 가 $t = 0$ 에서 극대가 되고

함수 $h(t) = f(a_1 + tm_1, \dots, a_n + tm_n)$ 가 $t = 0$ 에서 극소가 된다

는 것이 성립한다는 것이다.

기울기 = 0

^ 특징



미분가능한 다변수실함수의 정류점(기울기 벡터가 영벡터가 되는 점. 즉, 접평면이 수평이 되는 점)은, 안장점이 극값이다.