


Optimizer (옵티마이저)

Stochastic gradient descent. (SGD)
(확률적 경사 하강법)

$$W = W - \alpha \cdot \frac{dE}{dW}$$

weight learning rate gradient

$$\theta_{t+1} = \theta_t - \eta \cdot \nabla_{\theta} J(\theta)$$

Momentum

(모멘텀)

$$W = W + v - \alpha \cdot \frac{dE}{dW}$$

(0.9) velocity 이전 힘을 받지 않을 때 서서히
조정시키는 역할

AdaGrad

$$h = h + \frac{dE}{dW} \odot \frac{dE}{dW}$$

원소별 곱셈.

$$W = W - \alpha \frac{1}{\sqrt{h}} \cdot \frac{dE}{dW}$$

기울기의 제곱
매개변수 W 중에
가장 많이 움직인 원소
학습률을 낮아지게.

- 학습률 η 를 정하기 위한 방법
→ η 를 천천히 낮추다 모두에게
→ 모두보다 개별 맞춤으로
- 과거의 기울기를 저장해서 태워 때문에
어느 순간 경신량이 0에 가까워지는 현상
(RMS prop의 해결.)

Adam (RMS prop + Momentum)

- 2015년
- 하이퍼 파라미터의 편향 보정
- 지수평균 이동

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

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

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

$$\epsilon = 10^{-8}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_{t+1} + \epsilon}} \cdot \hat{m}_t$$

Adam

과거값을 줄여준다. How? \Rightarrow 지수 이동 평균을 통해!

초기 값의 속도가 작은 경우 ($\beta \approx 1$) $m(0), v(0) = 0$ 으로 초기화 되는데

이때 이동평균을 갱신하면 0으로 편향된다

값이 너무 작아 보정함

(β_1^t, β_2^t 는 각각 1번 곱함)
 $t \gg 1$ 일때 $\hat{v}^t = \frac{v^t}{1-\beta_1^t} = \frac{v^t}{1-\beta}$
 $(0.9)^{100} \approx 0$

Momentum + RMS

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

adagrad + RMS

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

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1-\beta_2^t}$$

$$\theta_{t+1} = \theta_t - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$$

$L=512$ step
 $\alpha = 0.001$
 $\beta_1 = 0.9$ = momentum decay rate
 $\beta_2 = 0.999$ = adaptive term decay rate
 $\epsilon = 10^{-8}$

기존의 momentum 과 달리 지수이동평균을 사용하게 ($1-\beta$)가 작다 작은 값을 갱신

기존 Momentum 과 달리

Rmsprop의 지수이동평균 사용

$\rightarrow (1-\beta)$ 이 작은 값을 갱신

Adam optimizer 보다 update 속도를 보정(올라가) 위해 다들 같이 잡음

$$\rightarrow \hat{\alpha}_t = \alpha \cdot \frac{\sqrt{1-\beta_2^t}}{\sqrt{1-\beta_1^t}}$$

이때 $n \gg 1$ 일때 $\hat{\alpha}_t \approx \alpha$ 로 보임

ex) $n = \infty$

$$\hat{\alpha}_t = \alpha \cdot \frac{\sqrt{1-\beta_2^t}}{\sqrt{1-\beta_1^t}} = \alpha \cdot 1$$

$$\therefore \hat{\alpha} = \alpha$$

learning rate //

ex) $f(x) = (2x-10)^2 \rightarrow W$ 환경에 대해서

init $x_0 = 0, lr = 0.1$ momentum decay rate $\beta_1 = 0.9$

adaptive term decay rate $\beta_2 = 0.999$

$$\nabla f(x) = 2(2x-10) \cdot 2 = 4(2x-10)$$

$$\bigcirc a_1 = (0.9) \cdot 0 + (1-0.9) \cdot 4(2x_0-10), (x_0=0), (a_0=0) = -4$$

$$\bigcirc b_1 = (0.999) \cdot 0 + (1-0.999) \cdot 16(2x_0-10)^2, (x_0=0) (b_0=0) = 0.001 \cdot 1600 = 1.6$$

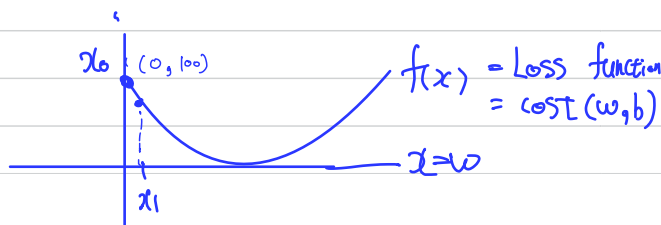
$$\bigcirc \hat{\alpha}_1 = (0.1) \cdot \frac{\sqrt{1-(0.999)^1}}{\sqrt{1-(0.9)^1}}, \alpha = 0.1$$

$$= 0.0158 \cdot \left(\frac{\sqrt{0.001}}{\sqrt{0.1}} \right) = \sqrt{\frac{0.001}{0.1}} = \sqrt{\frac{1}{100}} = \frac{1}{10} = 0.1$$

$$\bigcirc x_1 = x_0 - \frac{\hat{\alpha}_1}{\sqrt{b_1 + \epsilon}} \odot a_1, \epsilon = 1e-8$$

$$= 0 - \frac{0.0158}{\sqrt{1.6}} \cdot (-4)$$

$$= \frac{0.0632}{\sqrt{1.6}} = 0.0499 \dots$$



초기 update 속도가 GD 보다 느리지만 global minimum을 찾아지 않고 잘 도착함.

Momentum

weight를 업데이트 해줄때 이전 것들도 포함해 계산

- gradient가 0인 부분이 와도 이전 gradient가 남아있기에 움직이지.

$$V_t = \gamma \cdot V_{t-1} + \eta \cdot \nabla_{\theta_t} J(\theta_t)$$

$$\theta_{t+1} = \theta_t - V_t$$

이전 w들을 저장

$$\theta_{t+1} = \theta_t - \underbrace{\gamma \cdot V_{t-1}}_{\frac{\partial J}{\partial \theta}} - \underbrace{\eta \cdot \nabla_{\theta_t} J(\theta_t)}_{\frac{\partial J}{\partial \theta}}$$

γ = momentum decay rate
(≈ 0.9)

아 그래서 부드럽게 가는 거지
stochastic gradient descent는
그래서 w가 옳은 방향으로 가지만
그런데 확률적으로 적의가 강하면
지그재그로 가지!



ex) $\eta \cdot \nabla_{\theta_t} J(\theta_t) = f_t$

$$V_1 = f_1$$

$$V_2 = \gamma \cdot f_1 + f_2$$

$$V_3 = \gamma^2 \cdot f_1 + \gamma \cdot f_2 + f_3$$

$$V_4 = \gamma^3 \cdot f_1 + \gamma^2 \cdot f_2 + \gamma \cdot f_3 + f_4$$

($\gamma = 0.9$
= momentum decay rate)

→ 오래된 gradient의 영향력이 줄어든다.

- gradient가 0인, local minimum에서 벗어나
loss function을 더 탐험할 수 있도록 해준다.

Adagrad

- Momentum, SGD와 같은 optimizer는 η 이 고정하기에 global minimum을 지나칠 수도 있다.

→ 모든 파라미터에 대해 같은 learning rate를 천천히 낮춘다 (η)

→ 이용이되면 sparse한 node가 있을수도 있기에 개별 맞춤으로 (η)를 낮춘다.

ex) i 번째 input node 값을 $a_1, a_2, a_3, \dots, a_n$ 이라고 할때
 $i+1$ 번째 layer의 입력값은

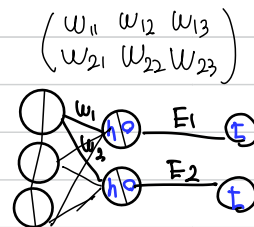
$$\underbrace{w_0}_{\text{bias}} + \underbrace{w_1 a_1}_{\text{weight}} + \underbrace{w_2 a_2}_{\text{input}} + \dots + w_n a_n \quad (w_i = \theta_i)$$

이때 만약 a_2 는 sparse해서 주로 0값이 많이 나올때.

학습할 데이터의 특징이 그 부분만 0값일때

$$w_0 + w_1 a_1 + 0 + \dots + w_n a_n$$

이때 w_2 term이 없어서 w_2 는 update 되지않음!



그러다 다음 데이터에서 오랜만에 $a_2 \neq 0$ 인 값이 나오면

이미 w_2 제외 w_n 들은 global minimum에 도달했지만
 w_2 는 거의 처음 위치에서 update 시작!

$$\frac{dE}{dw_2} = \frac{dE}{dh} \frac{dh}{dw_2} \quad , \quad (h = w_0 + w_1 a_1 + 0 + w_n a_n)$$

$$\frac{d}{dw_2} (w_0 + w_1 a_1 + 0 + w_n a_n) = 0$$

따라서 w_2 는 다른 파라미터 (w_n)를 더 크게 update 해주어야 global minimum에 조금이라도 빨리 들어갑니다!

- t 번째 step의 i 번째 파라미터 $\theta_{t,i}$
 (epoch, iterate) (weight) ($w_{t,i}$)

- t 시점의 $\theta_{t,i}$ 에 대한 gradient vector

$$\nabla_{\theta} J(\theta_{t,i}) = g_{t,i}$$

$$\frac{d \text{Loss}(w_{t,i})}{dw_{t,i}} = g_{t,i} \quad \Rightarrow \quad \text{Loss 함수의 } w_{t,i} \text{에 대한 기울기}$$

$$\text{SGD} \quad \theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}$$

$$\text{Adagrad} \quad \theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,i,i} + \epsilon}}$$

$$\text{Adagrad vector form} \quad \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

(η : learning rate)

(ϵ : $1e-8$)

($G_{t,i,i}$: $g_{1,i}^2 + g_{2,i}^2 + g_{3,i}^2 + \dots + g_{n,i}^2$)

Sparse 하다.

⇒ 특정 공간 안에 굉장히 상대적으로 적은 양의 값만 존재하는 경우 (드문)(희박한)

$$h = h + \frac{dE}{dw} \odot \frac{dE}{dw}$$

원소별 곱셈.

기울기의 제곱

$$w = w - \alpha \frac{1}{\sqrt{n}} \cdot \frac{dE}{dw}$$

매개변수 w 중에 가장 많이 움직인 원소 학습률을 낮아지게.

Rms Prop (Root Mean Square Propagation)

Adagrad는 η 가 증가하면서 $G_{t,i}$ 이 커져 learning rate가 소멸되는 문제점
(step, epoch, iterate) $\left(\frac{\partial L}{\partial w} \odot \frac{\partial L}{\partial w}, (\nabla_{w_i} J(\theta))^2 \right)$ (G_t 가 분모에 있기 때문)

↓ 나거 모든 기울기를 균일하게 더하지 않고 새로운 기울기의 정보만 반영하도록 해서 학습률이 0에 가까워지는 것을 방지

i) $G_{t,i}$ 를 저장하지 않고 **지수평준화**를 저장 (저장의 함)

ii) 과거 gradient의 **영향력**을 감소시키기 위해 **지수평준화**를 사용.

$$\chi_k = \alpha p_k + (1-\alpha)\chi_{k-1} \quad \text{where } \alpha = \frac{2}{N+1}$$

χ = 이동 평균 값

p = 현재 값 (w ?)

α = 영향력 (w에 대한 가중치?)

k = step, epoch

N = 값의 개수

ex) 처음부터 지금까지 계산 해보자

$$\begin{aligned} \chi_k &= \alpha \cdot p_k + (1-\alpha)\chi_{k-1} \\ &= \alpha \cdot p_k + (1-\alpha)(\alpha p_{k-1} + (1-\alpha)\chi_{k-2}) \\ &= \alpha \cdot p_k + \alpha(1-\alpha)p_{k-1} + (1-\alpha)^2\chi_{k-2} \\ &= \alpha \cdot p_k + \alpha(1-\alpha)p_{k-1} + (1-\alpha)^2(\alpha p_{k-2} + (1-\alpha)\chi_{k-3}) \\ &= \alpha \cdot p_k + \alpha(1-\alpha)p_{k-1} + \alpha(1-\alpha)^2p_{k-2} + (1-\alpha)^3\chi_{k-3} \\ &= \dots \\ &= \alpha(p_k + (1-\alpha)p_{k-1} + (1-\alpha)^2p_{k-2} + (1-\alpha)^3p_{k-3} + \dots + (1-\alpha)^{k-1}p_1) + (1-\alpha)^k\chi_{k-N} \\ &= \alpha \cdot \sum_{i=0}^{k-1} (1-\alpha)^i \cdot p_{k-i} + (1-\alpha)^k \cdot \chi_{k-N} \end{aligned}$$

이때 $(1-\alpha)$ 이 주피마과 곱해지고 $(1-\alpha) < 1$ 이기에 시간이 지나면 영향력이 작아진다.

저장된 이전 weight 들이

필터에문해 이런 α (가중치 영향력)

forgetting factor or decaying factor 라고 불림

$$\begin{aligned} \text{Rms prop} \quad \eta_t &= \eta \cdot G_{t-1} + (1-\eta) \left(\nabla_{w_i} J(\theta_t) \right)^2 \\ w_{t+1} &= w(t) - \eta \frac{1}{\sqrt{G(t) + \epsilon}} \cdot \nabla_{w_i} J(\theta_t) \end{aligned}$$

$$\begin{aligned} \left(\nabla_{w_i} J(\theta_t) \right)^2 &= \frac{\partial \text{Error}}{\partial w_{ti}} = \frac{\partial \text{Loss}(w_{ti})}{\partial w_{ti}} \\ \eta &= \alpha = \text{영향력 (현재 기울기)} \end{aligned}$$

9 이전 변화량과 현재와 변화량의 **어떠한 평균**으로 learning rate가 급격히 감소되는 걸 막을 수 있다

9 G_t 로 인해 최근 변화량에 **높은 영향력**을 준다 (r)

ex) $\eta = 0.9$

$$\begin{aligned} g_{ij}^{(t)} &= (0.9) \cdot g_{ij}^{(t-1)} + (0.1) \left(\frac{\partial L}{\partial w_{ij}} \right)^2 \\ &= (0.9) \left((0.9) \cdot g_{ij}^{(t-2)} + (0.1) \left(\frac{\partial L}{\partial w_{ij}} \right)^2 \right) + 0.1 \left(\frac{\partial L}{\partial w_{ij}} \right)^2 \\ &= (0.81) g_{ij}^{(t-2)} + (0.09) \left(\frac{\partial L}{\partial w_{ij}} \right)^2 + (0.1) \left(\frac{\partial L}{\partial w_{ij}} \right)^2 \end{aligned}$$

전전 값 (0.81) 전 값 (0.09) 현재