

09 & 10 . 경사 하강법

(Gradient Descent)

- Momentum
 - Netsterov Momentum
 - Ada grad
 - RMS Prop / Adadelta / Adam
- 학습률 고정
- 경사도 다른 변수의 학습률을 다르게!
- Momentum + Adagrad

경사 하강법 (Gradient Descent)

- Momentum

• 1st : 하강 방향 \ominus . α 는 학습률 (매번 계산되는 단계값을) 학습률로 고정!

$$v^{(1)} = -\alpha g^{(1)}, \text{ 속도 (velocity) 이동 방향 계산}$$

$$x^{(1)} = x^{(0)} + v^{(1)} = x^{(0)} - \alpha g^{(0)}, \text{ 경사 하강 (Gradient Descent)}$$

→ 속도 : 최초의 계획점이 주어지면 거기서의 경사도 계산 후, 주어진 학습률에 의해 하강 방향으로 이동하는 값의 크기

• 2nd :

$$v^{(2)} = \beta v^{(1)} - \alpha g^{(2)} = \beta (\alpha g^{(1)}) - \alpha g^{(2)}, \text{ 가속도 (Momentum)}$$

$$x^{(2)} = x^{(1)} + v^{(2)} \text{ Momentum 파라미터}$$

• Kth :

$$v^{(k+1)} = \beta v^{(k)} - \alpha g^{(k)} = \beta (-\alpha g^{(k-1)}) - \alpha g^{(k)}$$

$$x^{(k+1)} = x^{(k)} + v^{(k+1)} \rightarrow \text{Momentum update formula}$$

→ 최소값이 있는 방향으로 계획점 꾸준히 update

But, 빨리, 큰 step으로 이동 & 최소값 안정적 계산!

⇒ Momentum 방법은 이전 반복에서 얻어진 경사 하강 방향으로의

진행 크기 (Velocity)를 누적!

최초에는 경사 하강법 (Gradient Descent)와 똑같이 계획점 이동

2nd는 이전에 얻어진 Velocity를 새로 얻어진 Velocity에 누적

완전한 목적함수 경우

→ 좀 더 빨리 가고 싶는데 이전에 진행했던 방향 크기 누적해 진행!

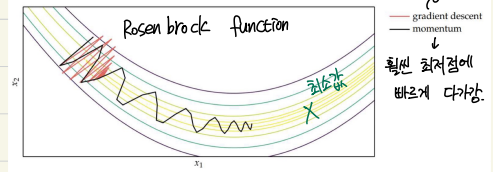
⇒ 경험적으로, β 는 1일 때 가속도 커져 1보다 작게! 0.9 씩 이용

따라서, 실제 hyper-parameter은 학습률 (α)

$$v^{(k+1)} = \beta v^{(k)} - \alpha g^{(k)}, \beta = 0.9$$

$$x^{(k+1)} = x^{(k)} + v^{(k+1)}$$

Gradient Descent는 사이클링 발생.



→ long deep valley 있으면 Momentum 방법 좋음

다른 함수에서는 Gradient Descent와 Momentum 비슷

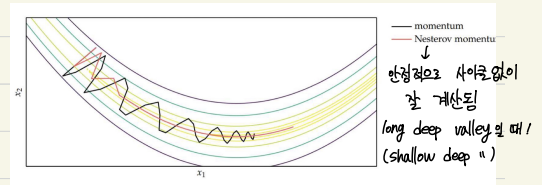
- Nesterov Momentum

$$v^{(k+1)} = \beta v^{(k)} - \alpha g^{(k)} = \beta v^{(k)} - \alpha \nabla f(x^{(k)})$$

$$\begin{cases} v^{(k+1)} = \beta v^{(k)} - \alpha \cdot \nabla f(x^{(k)} - \beta v^{(k)}) \\ x^{(k+1)} = x^{(k)} + v^{(k+1)} \end{cases}$$

Nesterov Momentum이 완전한 경사 방향으로

큰 step size 보이게 됨



gradient

⇒ Gradient Descent, Conjugate Descent,

Momentum, Nesterov Momentum 모두

다변량 변수에 각각 동일한 학습률 부여

변수가 커지며 변화가 큰 변수의 경사도 방향으로 골려다님

→ 변수 별로 다른 가중치 주며 계산값 끌어다여 함.

경사 하강법 (Gradient Descent)

- Adagrad : Adaptive subgradient (2011)

- χ_i 의 크기가 다르므로 χ_i 에 서로 다른 α

$$\bullet \chi_i^{(k+1)} = \chi_i^{(k)} - \left(\frac{\alpha}{\sqrt{S_i^{(k)}} + \epsilon} \right) g_i^{(k)}$$

$S_i^{(k)}$ 에 따라 변수별 학습률이 다르게 부여

$$\rightarrow S_i^{(k)} = \sum_{j=1}^k (g_i^{(j)})^2, \quad 1\text{번째 변수의 } k\text{번째까지의 제곱합}$$

$$\rightarrow \epsilon = 1 \times 10^{-8} = 1E-8 \quad \text{각 변수별 수직}$$

↳ 분모가 0이 되는 것 방지

\rightarrow 변수들이 gradient 양이 크면 $S_i^{(k)}$ 커져나가 학습률 작아짐

- k 가 증가할수록 $\frac{\alpha}{\sqrt{S_i^{(k)}}}$ 가 지속적으로 작아짐 (decaying)

\rightarrow 업데이트 안되는 문제! (Vanishing)

- RMS Prop : Geoffrey Hinton

- 경사도 제곱 크기 유지 \hookrightarrow gradient 제곱되는 dot product

$$\bullet \hat{S}^{(k+1)} = \gamma \hat{S}^{(k)} + (1-\gamma) (g^{(k)} \odot g^{(k)}), \quad \gamma = 0.9$$

$\hookrightarrow \gamma$ 와 $1-\gamma$ 로 해서 합=1 $\hat{S}^{(k+1)}$ 커지는 것 막음

$$\beta, \alpha \rightarrow \gamma, (1-\gamma) \quad \text{gradient의 가중치 줄임 (가중평균)}$$

$$\chi_i^{(k+1)} = \chi_i^{(k)} - \frac{\alpha}{\sqrt{S_i^{(k)}} + \epsilon} \cdot g_i^{(k)}$$

$$= \chi_i^{(k)} - \frac{\alpha}{\text{RMS}(g_i) + \epsilon} \cdot g_i^{(k)}$$

$$\rightarrow \text{RMS}(g_i) \text{ 커지면 } \frac{\alpha}{\text{RMS}(g_i) + \epsilon} \text{ 작아진 위험 있음}$$

- Adadelta : Adaptive Learning Rate Methods

- 학습률 : RMS 증가 비율

$$\bullet \hat{S}^{(k+1)} = \gamma \hat{S}^{(k)} + (1-\gamma) (g^{(k)} \odot g^{(k)}), \quad \gamma = 0.9$$

$$\chi_i^{(k+1)} = \chi_i^{(k)} - \frac{\text{RMS}(\Delta \chi_i)}{\text{RMS}(g_i) + \epsilon} g_i^{(k)}$$

$\rightarrow \text{RMS}(\Delta \chi_i)$ 는 $\hat{S}^{(k)}$ 과 $\hat{S}^{(k+1)}$ 의 차이 제곱

$\frac{\alpha}{\sqrt{S_i^{(k)}}}$ 가 작아지는 문제 완화.

변수마다 동일 α 적용되는 문제 완화

- Adam : Adaptive momentum estimation method

$$\bullet \text{momentum} : \nu^{(k+1)} = \gamma_\nu \nu^{(k)} + (1-\gamma_\nu) g^{(k)}$$

$$\bullet \text{Squared gradient} : \hat{S}^{(k+1)} = \gamma_s \hat{S}^{(k)} + (1-\gamma_s) (g^{(k)} \odot g^{(k)})$$

$\rightarrow \nu^{(k+1)}, \hat{S}^{(k+1)}$ 추정식 사용하자!

$$\bullet \text{momentum correction} : \hat{\nu}^{(k+1)} = \frac{\nu^{(k+1)}}{(1-\gamma_\nu)^k} \quad (\text{like 평균})$$

\hookrightarrow momentum 이용. k 증가 \rightarrow 분모는 1과 가까운 값

$$\bullet \text{Squared gradient correction} : \hat{\hat{S}}^{(k+1)} = \frac{\hat{S}^{(k+1)}}{(1-\gamma_s)^k}$$

\hookrightarrow 값이 커진 위험이 있어서 전체 반복에서 동일한 평균값 추정
like 분산

$$\bullet \text{Update} : \chi^{(k+1)} = \chi^{(k)} - \frac{\alpha \hat{\nu}^{(k+1)}}{\sqrt{\hat{\hat{S}}^{(k+1)}} + \epsilon}$$

\downarrow

계산량 \nearrow . 안정적으로 빠르게 최소값에 접근해

Computational Cost 줄어둘 수 있음