

딥러닝 Optimization 방법론

↳ Neural Net을 학습시키는 training 방법.

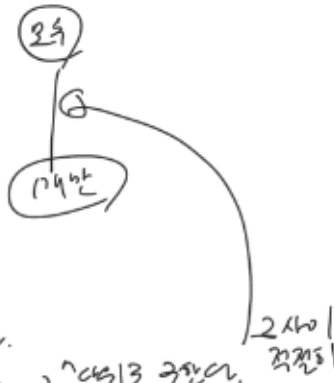
① Gradient Descent \rightarrow 어떤 함수 U 에 대해 최소값 찾기

↳ 내 parameter의 기울기를 알고 cost를 최소화하는 기울기를 찾자

Batch gradient descent
 \Rightarrow 모든 data 가용하고 평균

Stochastic
 \Rightarrow 한번에 하나만 가용하기

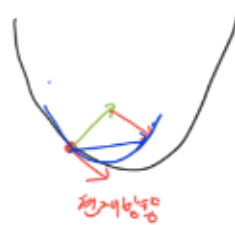
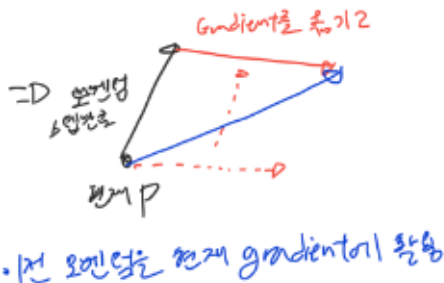
Mini-batch
 \Rightarrow 중간값점들의 가용성을 구한다.
 ex) 50 ~ 250개 정도 2^k 형태로 구한다.



* 적절한 learning rate을 정한다.

가용데이터

* ① Momentum \Rightarrow 이전 momentum + 현재 Gradient를 다음 Gradient에 반영.



② Nesterov accelerated gradient

exponential moving average 을 쓰지않고

* Adagrad = learning rate 스케줄링이 \Rightarrow 모든 parameter에 learning rates를 바꾼다.

para $\left(\begin{matrix} 1000000 \\ 177 \end{matrix} \right)$ \Rightarrow 많이변한 para는 조금 변화하기
 조금 " " 많이 "

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,i}}} g_{t,i}$$

\Rightarrow Gradient의 제곱 \Rightarrow G가 크면 $\theta_{t+1,i}$ 는 작아질 (분모바라)

* 문제점 G가 계속 커지면 $\theta_{t+1,i}$ 가 계속 작아질
 \Rightarrow 학습 더 이상 안됨.

||

→ 84/12

$$\begin{cases} E[g^2]_t = \sigma E[g^2]_{t-1} + (1-\sigma)g_t^2 \\ E[\Delta\theta^2]_t = \sigma E[\Delta\theta^2]_{t-1} + (1-\sigma)\Delta\theta_t^2 \end{cases}$$

$$\theta_{t+1} = \theta_t - \frac{\sqrt{E[\Delta\theta^2]_t + \epsilon}}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

Adadelta

learning rate

가 없는 버전

임실론 문제: 초기값이 0가 이므로
epsilon 주면 처음 다음 같은 설정

* RMSprop

learning rate

가 있는 버전

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

global learning rate

⑧

Adagrad: Exponential moving average 'n'을 사용하지 않고 Gradient 제곱으로 θ_{t+1} 조정
⇒ adaptive learning rate를 가능하지만 learning rate = 1 계속 감소

↓ ①

RMSprop: Exponential moving average를 사용해서
⇒ global lr을 가짐

이를 막기위해

adadelta에서

exponential moving avg를

사용하여 최근 Gradient가

변화하는 것을 check

RMSprop과 Adadelta의 차이 ②

⇒ global learning rate = 1이라
"n"에 따라

* Adam ⇒ learning rate와 momentum을 합침

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

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

g의 제곱

g의 제곱

adaptive lr을 위해 따라서 G를 expo avg한 것

보통 0.9이상의 사용

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t + \epsilon}} \times \frac{\sqrt{1-\beta_2^t}}{1-\beta_1^t} m_t$$

동적화강

unbias based를 보정하기위해서
0이 가까워질

임실론 → 초기 g가 0 이되는 것을 막기위한 값

adam에서 이 초기값이 e^{-8} 정도인데 e^{-4} 정도가 학습이 잘됨 (첫제곱)

⑧

반반한 경우에는

Adam을 쓸 때 optimize하면

오버피팅은 사용 (이전 기호기 사용) gradient는

↓
 m 사용하기 싫으면 RMSprop ⇒ global exponential moving avg
 사용 = 'g' 제곱

질문 : optimizer에 어떤 방법을 쓰는 것이 좋은가?

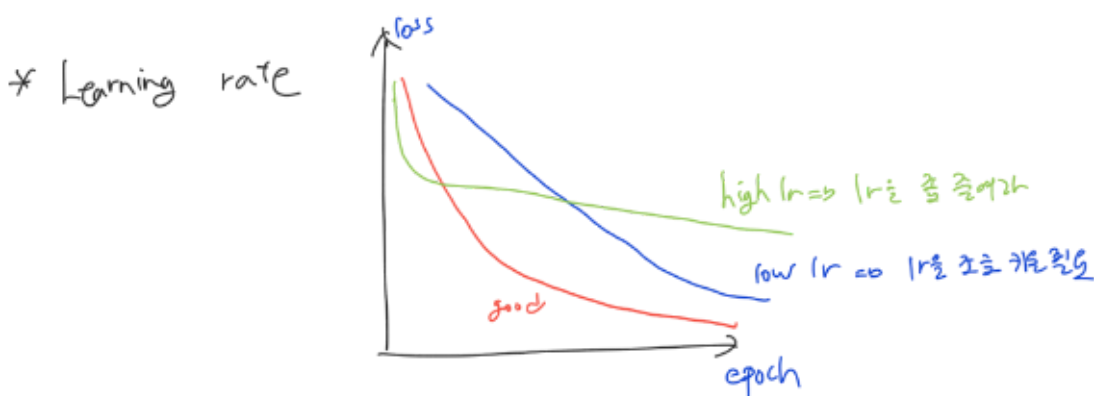
⇒ Adaptive learning rate를 사용하는 것이 좋다.

(Adagrad, Adadelta, RMSprop, Adam)

Adagrad ⇒ learning rate이 계속 감소하는 특징 why? G (기호기제곱)이 너무 크

RMSprop ⇒ Adagrad의 lr 감소 해결 (Adadelta로 바꿀거지)

Adam은 bias-correction과 momentum을 RMSprop에 추가한 것
 이전 기호기 사용 $\frac{\sqrt{1-\beta^t}}{1-\beta}$ \Rightarrow bias 보정
 β ⇒ 오버피팅 구하기 사용



* Restricted Boltzmann Machine (RBM)

① Energy-based models = unsupervised model (img만 주고 어떻게 관련관지
 관련 모델)

$$p(x) = \frac{e^{-E(x)}}{Z}$$

↓
 기호기 입력, $E(x)$ 는 입력이 주어졌을 확률이고
 e : energy의 $E(x)$ 는 반비례이므로 '-'
 z : normal parameter

* Supervised learning이란
 (지도학습)

입력 \rightarrow $\boxed{+}$ \rightarrow 출력

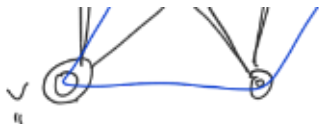
hidden

h



∴ 목적은 E 를 학습하는 것.

\boxed{h} hidden layer



visible layer

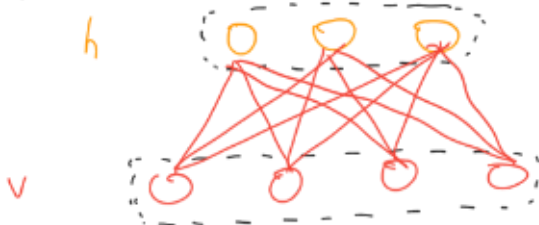
visible = 우리가 갖고 있는 Image

Restricted Boltzmann Machine은

= Restricts connection between visible and hidden

즉, RBM은 $\left(\begin{array}{l} v \text{ 개리 연결 } \times \\ h \text{ 개리 연결 } \times \end{array} \right)$ v 와 h 만 서로 연결

RBM



* Energy (RBM에서는 에너지를 학습하면 됨)

$$E(v, h | \theta) = - \sum_{i=1}^D \sum_{j=1}^F w_{ij} v_i h_j - \sum_{i=1}^D b_i v_i - \sum_{j=1}^F a_j h_j$$

w_{ij}, b_i, a_j 를 구하는 것.

* Probability (확률)

$$P(v, h | \theta) = \frac{1}{Z(\theta)} \exp(-E(v, h | \theta))$$

normalization constant $Z(\theta)$

$$Z(\theta) = \sum_{v, h} \exp(-E(v, h | \theta))$$

(모든 가능한 v, h layer에 대한 에너지의 exp한 것)

$P(v | y) \Rightarrow y$ 가 주어졌을 때 v 가 나올 확률

$$= \frac{P(v, y)}{\sum_x P(x, y)}$$

즉 x 가 y 일 때 y 의 확률 $= P(y)$

* v 가 주어졌을 때 v 를 최대화하는 parameter 찾기

$$w_{ij}, b_i, a_j$$

$$\hat{\theta} = \operatorname{argmax} \log P(v | \theta)$$

$$\text{FO} \frac{\partial}{\partial \theta} (-\log P(v)) = \frac{\partial}{\partial \theta} (-\log \sum_h P(v, h))$$

sum rule

$$\frac{d}{d\theta} \log f(\theta) = \frac{1}{f(\theta)} \cdot \frac{d}{d\theta} f(\theta)$$

$$= \frac{\partial}{\partial \theta} \left(-\log \sum_h \frac{1}{Z} \exp(-E(v, h)) \right)$$

즉 θ 의 함수, E 도 θ 의 함수이므로 θ 에 대한 함수를 미분

$$\frac{d}{d\theta} (f(\theta)g(\theta)) = f'(\theta)g(\theta) + f(\theta)g'(\theta)$$

$$= - \frac{1}{Z} \left(\sum_h \frac{\partial}{\partial \theta} \exp(-E(v, h)) - \sum_h \frac{1}{Z^2} \frac{\partial Z}{\partial \theta} \exp(-E(v, h)) \right)$$

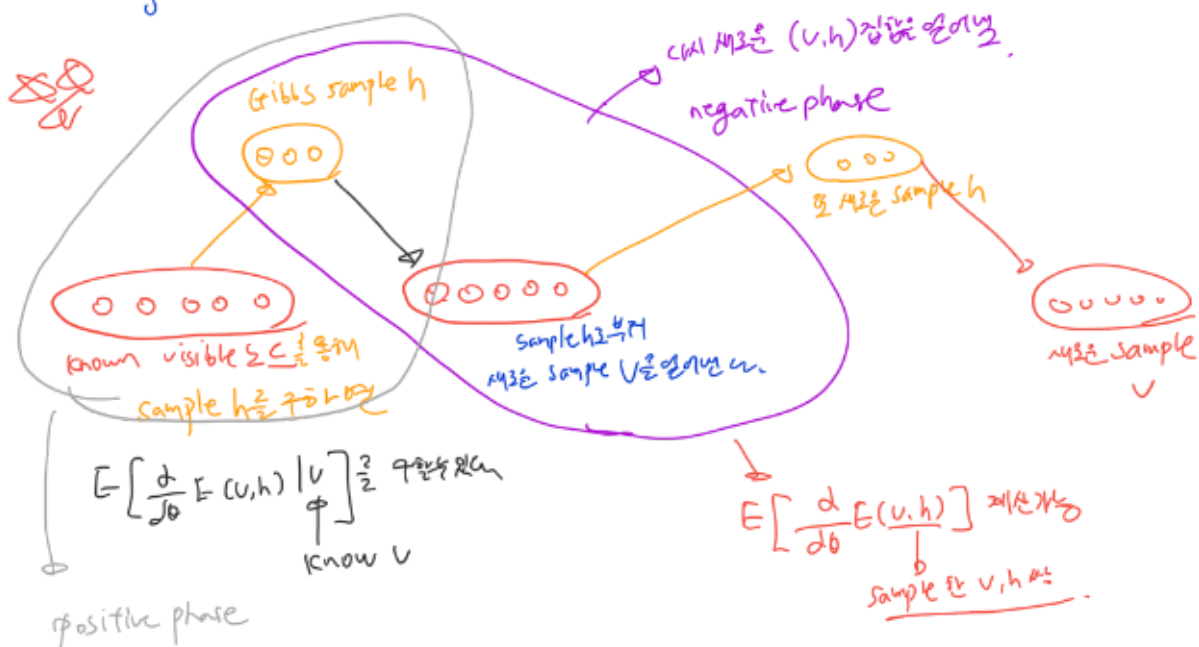
$\frac{\partial}{\partial \theta} (-\log p(v)) = \sum_h \left(\frac{\exp(-E(v,h))}{\sum_{h'} \exp(-E(v,h'))} \cdot \frac{\partial}{\partial \theta} E(v,h) \right) + \frac{1}{Z} \frac{\partial Z}{\partial \theta}$

$Z(\theta) = \sum_{v,h} \exp(-E(v,h))$

$p(x,y) = \frac{p(x,y)}{\sum_x p(x,y)}$

$p(x|y)$ 이므로

$$= E \left[\frac{d}{d\theta} E(v, h) \mid v \right] - E \left[\frac{d}{d\theta} E(v, h) \right]$$



많은 것들이
줄곧 오를 때 학습시켜서
점점과 비슷하게
반응을 학습

마지막 수정: 오후 5:56