



## # 4. Optimization : most important step in training

→ loss  $\hat{y}$  & correct parameter  $\theta$ 를 통해 학습하는 과정

### ## Introduction

#### 1. Optimization in deep learning

→ training 데이터 optimization

→ training error, object function을 최소화하는 문제

- most difficult optimization task in DL : training

→ important and expensive

- main stream : stochastic gradient descent and its variants
  - + second-order methods, convex optimization

#### 2. Derivatives and optimization order

- derivatives

① first derivatives (= gradient)  $\Rightarrow$  slope

② second derivatives (= Hessian)  $\Rightarrow$  curvature

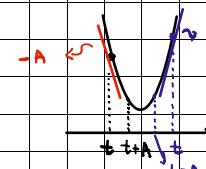
- optimization

① first-order algorithms  $\Rightarrow$  use only gradients  $\Rightarrow \mathbf{w} \leftarrow \mathbf{w} - \eta \nabla J$

② second-order algorithms  $\Rightarrow$  also use Hessian matrix  $\Rightarrow \mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{H}^{-1}$

#### 3. Critical points (= stationary points)

- points with zero slope :  $\nabla_x f(x) = 0$



→ derivative gives no info about which direction to move

→ three types : maxima (- curvature), minima (+ curvature), saddle point

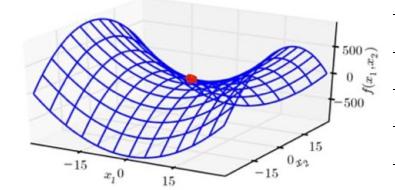
- a saddle point : contains both positive and negative curvature

$$f(\mathbf{x}) = x_1^2 - x_2^2 \quad \rightarrow \nabla f(2x_1, -2x_2) = \nabla f(0,0) = 0$$

$$\mathbf{H} = \nabla^2 f = \begin{pmatrix} 2 & 0 \\ 0 & -2 \end{pmatrix} \quad \rightarrow \text{eigenvalues} : \lambda_1 = +2, \lambda_2 = -2$$

(+), (-) 혼재!

saddle point!

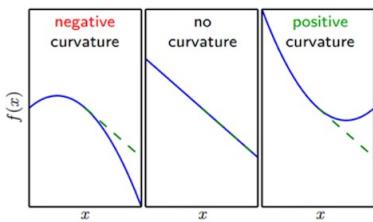


## 4. Use of second derivative

① to characterize critical points : min, max, saddle

② to measure curvature

③ to predict performance of an update in gradient based opti



- **negative curvature**
  - $f$  decreases faster than gradient predicts
- **no curvature**
  - gradient predicts the decrease correctly
- **positive curvature**
  - $f$  decreases slower than gradient predicts (eventually increases)

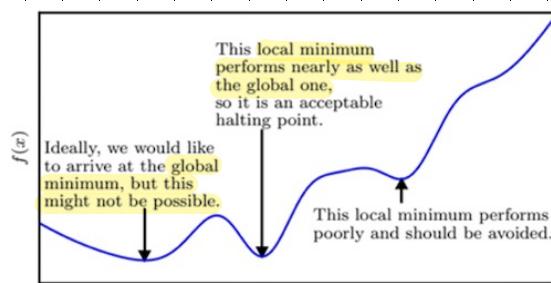
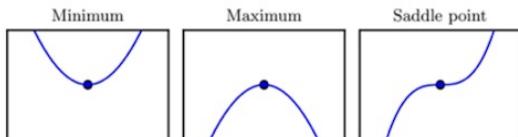
## 5. In deep learning

• Our objective function has

⇒ many local minima + many saddle points surrounded by very flat regions



⇒ makes optimization very difficult ex) → saddle point → gradient = 0 → training X



⇒ in high dim : Saddle points are much more common than local minima

## ## Gradient-Based Optimization

### Gradient Descent and its limitations.

## 6. Method of gradient descent.

- derivative : useful for minimizing a function (objective function)

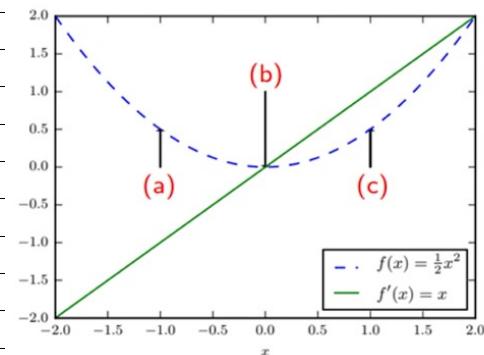
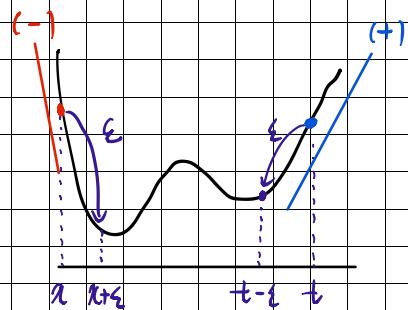
$$\Rightarrow \text{for small } \epsilon : f(x - \epsilon \cdot \text{sign}[f'(x)]) < f(x)$$

new para      old para

- We can thus reduce f(x) by

$\Rightarrow$  moving  $x$  in small steps with opposite sign of derivative.

= method of gradient descent.



$$\underbrace{x'}_{\text{new}} = \underbrace{x}_{\text{old}} - \underbrace{\epsilon}_{\text{learning rate}} \nabla_x f(x)$$

(a)  $x < 0: f'(x) < 0$   
 $\Rightarrow$  can decrease  $f$  by moving rightward

(b)  $x = 0: f'(x) = 0$   
 $\Rightarrow$  gd halts here (global min)

(c)  $x > 0: f'(x) > 0$   
 $\Rightarrow$  can decrease  $f$  by moving leftward

## 7. Sgd and its variants

- probably the most used optimization can obtain an unbiased estimate of gradient.

### Algorithm 1 gradient descent

```

1: initialize  $\theta$ 
2: while stopping criterion not met do
3:   sample  $m$  examples:  $\mathbb{X}_m = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ 
4:   compute gradient estimate:  $\hat{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)})$   $\triangleright m$  forward props
5:   apply update:  $\theta \leftarrow \theta - \epsilon \hat{g}$   $\epsilon$ : learning rate
6: end while

```

$$\begin{aligned} \nabla L(\theta) &= \frac{1}{N} \sum_i \nabla l_i(\theta) \\ g_B(\theta) &= \frac{1}{M} \sum_{i=1}^M \nabla l_i(\theta) \\ E_B[g_B(\theta)] &- E_B\left[\frac{1}{N} \sum_{i=1}^N \nabla l_i(\theta)\right] \\ &= \frac{1}{M} \sum_{i=1}^M E_B[\nabla l_i(\theta)] \\ &= \frac{1}{M} \cdot M \cdot \nabla L(\theta) \\ &= \nabla L(\theta)_B \end{aligned}$$

- tree variants ( $N$ : total # of examples)

①  $m=1$  : sgd

이들을 차례로 만날게.

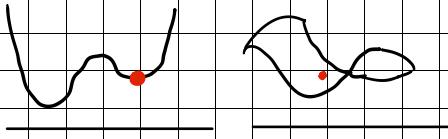
②  $1 < m < N$  : minibatch sgd

③  $N = m$  : batch gradient descent

## 8. Properties of SGD: bad ones

- SGD may suffer in the following situations

① local minima / saddle points

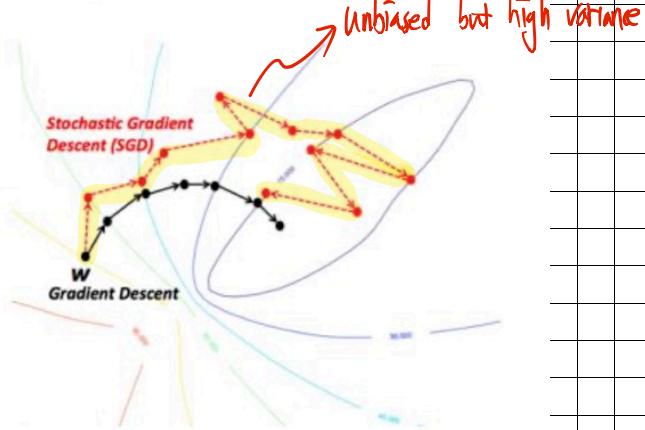


② gradient noise

because of minibatch ?? → poor condition  $\text{cond}(H) \neq \infty$

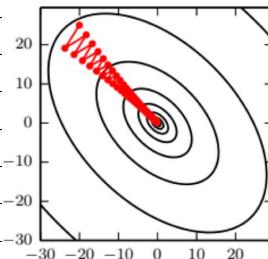
→ EWMA ... →

➤ gradient noise



③ Zero gradient  $\rightarrow$  gradient descent gets stuck

$\Rightarrow$  poor conditioning of  $H$



## 9. Poor conditioning of H

↳ 두 번 다른 주름 정보를 담고 있는 matrix

- Consider a point  $x$  in multiple dimensions : different second derivative for each direction.

- condition number of Hessian  $H$  at  $x$

⇒ measures how much the second derivatives differ from each other

$$\Rightarrow \text{condition number of } H \text{ at } X = \max_{ij} \frac{\lambda_i}{\lambda_j}$$

→ 가장 대  
→ 가장 소

- When  $H$  has a large condition number (poorly conditioned)

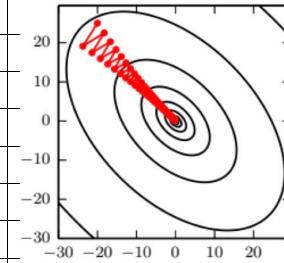
① gradient descent performs poorly

② it is difficult to choose a good step size  $\epsilon$

(e.g. N.M에서  $\frac{f'}{f''}$  term에서  $f''$ 를 구해도)

주름 정보를 실시간으로 반영

but 우리는 실시간으로 반영하는 주름을 fix한  $\epsilon$ 에 의지.



- how to handle poor conditioning without directly considering  $H$ ?

## Exponentially Weighted Average

### 10. Exponentially weighted moving average (EWMA)

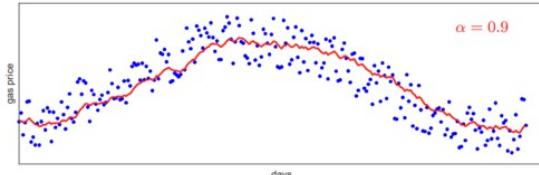
- given time series :  $g_1, g_2, \dots$

- EWMA defined as :

$$V_t = \begin{cases} g_1, & t=1 \\ \alpha V_{t-1} + (1-\alpha) g_t, & t>1 \end{cases}$$

$\alpha \in [0, 1]$  : degree of weighting decrease (constant Smoothing factor)

$$v_t = \alpha \cdot v_{t-1} + (1-\alpha) \cdot g_t$$



## II. Properties of EWMA

- effective weighting decreases exponentially over time.

→ 과거의 값의 영향력(?)이 차수별로 감소

$$\begin{aligned}
 V_t &= \alpha V_{t-1} + (1-\alpha) g_t \\
 k \curvearrowleft V_t &= \alpha [\alpha V_{t-2} + (1-\alpha) g_{t-1}] + (1-\alpha) g_t \\
 &\vdots \\
 &= \alpha^k V_{t-k} + (1-\alpha) [g_t + \alpha g_{t-1} + \cancel{\alpha^2 g_{t-2}} + \dots + \cancel{\alpha^{k-1} g_{t-k+1}}]
 \end{aligned}$$

→ Weight exponentially decreases toward the past

$\Rightarrow$  thus called "exponentially weighted"

- $$\bullet \text{ Recall } 1 + \alpha + \alpha^2 + \dots = \frac{1}{1 - \alpha}$$

$$\text{then } V_t = (1-\alpha) [g_t + \alpha g_{t+1} + \alpha^2 g_{t+2} + \alpha^3 g_{t+3} + \dots]$$

↳ effective number of observations

- higher  $\alpha$  (= more weight to past, less weight to present)

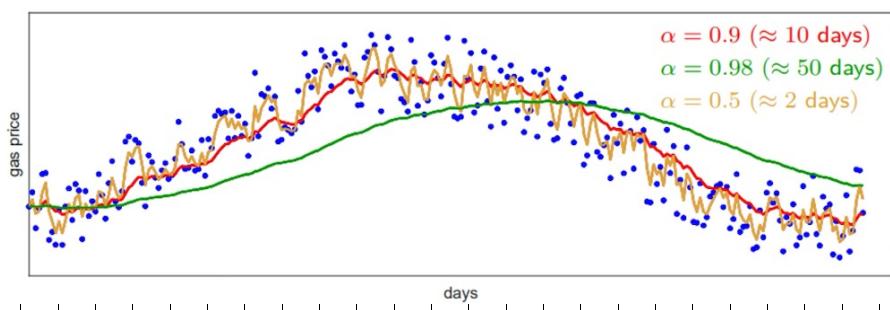
➤ smoother curve

← averaging over more days

➤ shifted further

← averaging over larger window

⇒ curve adapts more slowly to changes with more latency



## 12. Bias correction : Adam

⇒ 유의가 쓴 MA : EWMA

$$V_t = (1-\alpha)g_t + \alpha V_{t-1}, \quad (\alpha \text{ 를} \rightarrow 1 \text{에} \text{ 따라} \text{ 과거의} \text{ 영향} \text{을} \text{ 감소} \text{시} \text{킴})$$

$\frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}}$  학번. 2학년 학번 있는 것처럼  $\alpha$ 로 감소...  $\rightarrow 0$

$$\text{ex)} \quad V_t = (1-\alpha) \sum_{k=0}^{t-1} \alpha^k g_{t-k} \quad (\because \text{stationary} \Rightarrow E(g_t) = \mu)$$

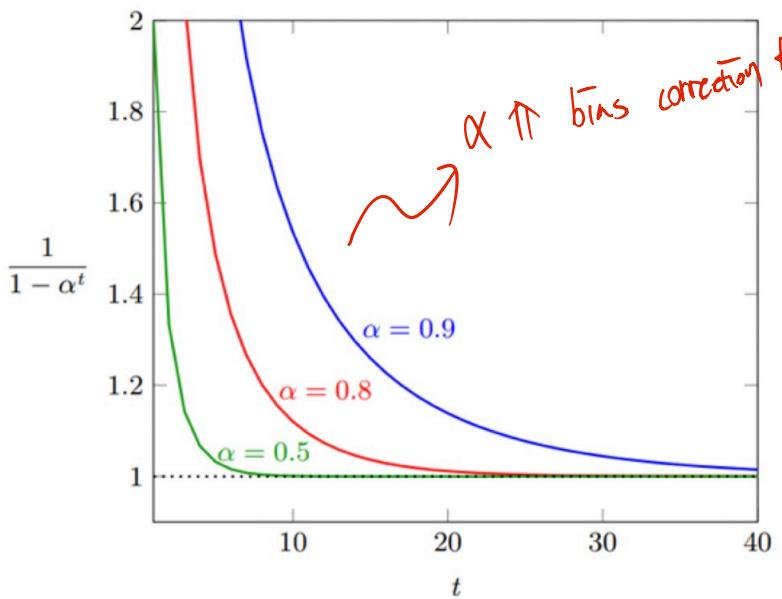
$$E(V_t) = (1-\alpha)\mu \sum_{k=0}^{t-1} \alpha^k \quad (\because \sum_{k=0}^{t-1} \alpha^k = 1 + \alpha + \dots + \alpha^{t-1} = \frac{1-\alpha^t}{1-\alpha})$$

$$= (1-\alpha^t)\mu$$

stationary  $\rightarrow$   $\mu \rightarrow$   $\text{Stationary} \rightarrow \mu$

$(1-\alpha^t)$   $\rightarrow$  biased  $\rightarrow$   $\text{unbiased}$

$$\text{纠正} \quad \hat{V}_t = \frac{V_t}{1-\alpha^t} \quad \text{이} \quad \text{임!}$$



## - Gradient Descent with Momentum

### 13. Method of momentum

- sgd: very popular but sometimes slow  $\Rightarrow$  unbiased but high variance

- Method of momentum  $\Rightarrow$  designed to accelerate learning  
especially in face of
  - ① high curvature  $\Rightarrow H_{\theta} \lambda_{\text{max}} \uparrow \Rightarrow$  poor cond
  - ② small but consistent gradients  $\Rightarrow$  saddle points
  - ③ noisy gradients

- Common algorithm

$\rightarrow$  accumulates exponentially decaying moving average of past gradients  
 $\Rightarrow$  gradient descent + grad. vgt. EWMA  
 $\rightarrow$  then continues to move in their direction.

### 14. Gradient descent with momentum

- Idea: compute EWMA of gradients and use it to update weights  
 $\rightarrow$  gradient  $\hat{g}$   $\hat{g}_t = \frac{\alpha}{1-\alpha} \hat{g}_{t-1} + \frac{1-\alpha}{1-\alpha} g_t$   $\xrightarrow{\text{grad. vgt.}} \text{EWMA of } \hat{g} \approx \frac{1-\alpha}{1-\alpha} g_t$   
 $\rightarrow$  noise  $\downarrow$ , consistent direction gradient.

- let  $g \triangleq \nabla_{\theta} J(\theta)$  ( $\theta$  represents  $W$  and  $b$  altogether)

$$\boxed{\text{form 1}} \quad \begin{aligned} v &\leftarrow \alpha v - \epsilon g \\ \theta &\leftarrow \theta + v \end{aligned} \quad \Rightarrow \quad \theta \leftarrow \theta - \epsilon \left( g + \frac{\alpha}{-\epsilon} v \right)$$

$$\boxed{\text{form 2}} \quad \begin{aligned} v &\leftarrow \alpha v + (1-\alpha)g \\ \theta &\leftarrow \theta - \epsilon v \end{aligned} \quad \xrightarrow{^3} \quad \theta \leftarrow \theta - \epsilon \left( g + \frac{\alpha}{1-\alpha} v \right)$$

$$\boxed{\text{form 3}} \quad \begin{aligned} v &\leftarrow \alpha v + g \\ \theta &\leftarrow \theta - \epsilon v \end{aligned} \quad \Rightarrow \quad \theta \leftarrow \theta - \epsilon (g + \alpha v)$$

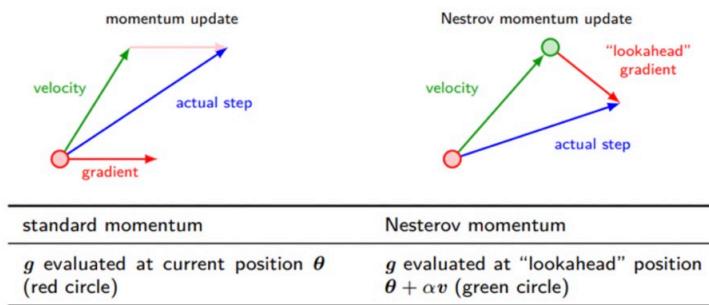
- bottom line:  $\theta$  is updated by linear combination of gradient and velocity

$$\theta \leftarrow \theta - \epsilon \left( \underbrace{g}_{\text{gradient}} + \text{constant} \cdot \underbrace{v}_{\text{velocity}} \right)$$

## 15. Nesterov momentum

- difference from standard momentum:

➤ `_where_gradient g = ∇θ J` is evaluated



- Nesterov momentum update rule

$$: V \leftarrow \alpha V - \varepsilon \nabla_{\theta} J(\theta + \alpha V)$$

$$\theta \leftarrow \theta + V$$

⇒ gradient is evaluated after the current velocity is applied

→ interpreted as adding a correction factor to standard momentum

momentum의 방향은?

gradient는 방향은?

↳ 45도 아래로 optim.SGD(momentum=0.9, nesterov=True)

$$V \leftarrow 0.9 V - \varepsilon \nabla_{\theta} J(\theta + 0.9 V)$$

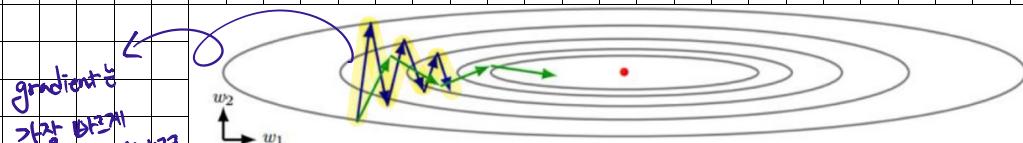
$$\theta \leftarrow \theta + V$$

## Pre parameter Adaptive Learning Rates : parameter $\theta$ & $\eta$ learning Rate $\epsilon$ !!

- optimization methods explained so far
  - set learning rate  $\epsilon$  globally and equally for all parameter :  $w \leftarrow w - \epsilon g$
- methods presented now : AdaGrad, RMSprop, Adam
  - adaptively tune  $\epsilon$  for each parameter

### 16. Motivation

- recall : limitation of gradient descent



gradient  
가장 빠른 방향  
증가 방향  $\nabla f(w)$   
& Minibatch?

→ goal : move horizontally

(cause by minibatch & curvature so high?)

→ problem : huge vertical oscillations

→ solution : we want to ① slow down learning vertically

② speed up learning horizontally

how?

without relying on  $H$  explicitly?

### 17. AdaGrad

- individually adapts learning rate of each direction per parameter

→ steep direction (large  $\frac{\partial J}{\partial \theta_j}$ ) : slow down learning

→ gently sloped direction (small  $\frac{\partial J}{\partial \theta_j}$ ) : speed up learning

이전에  $\sum_{t=1}^T \theta_j^{(t)}$   
이전에  $\sum_{t=1}^T \theta_j^{(t)} + \frac{1}{T} \sum_{t=1}^T \theta_j^{(t)}$   
 $\theta_j^{(t+1)} = \theta_j^{(t)} - \frac{1}{T} \sum_{t=1}^T \theta_j^{(t)}$   
not hyperparam

S: ① Adadelta: fix window

② EWMA

- adjusts learning rates per parameter

$$\epsilon_j = \frac{\epsilon}{\sqrt{\sum_{all \ past \ iter} (g_j \cdot g_j)}}$$

parameter which depends on past gradients ( $g^2$ )

learning rate of parameter  $\theta_j$

Hyperparameter  $\epsilon$

- net effect : greater progress in more gently sloped direction

p : monotonically decreasing : too aggressive  
⇒ stops learning too early!!

## 18. RMSProp (root-mean-square prop)

- modifies AdaGrad to perform better in non-convex setting
  - changes gradient accumulation to EWMA
  - mon  $\rightarrow g \text{ zt in EWMA}$
  - $\text{RMS} \Rightarrow g^2(\text{zt}) \text{ of EWMA}$
- use of exponentially decaying average allows RMSprop to
  - discard history from extreme past
  - ⇒ converge rapidly after finding a convex bowl

- Comparison ( $r$ : accumulation variable)

① AdaGrad

$$r \leftarrow r + g \otimes g$$

$$\Delta \theta_t \leftarrow -\frac{\epsilon}{\sqrt{s+r}} \otimes g$$

$$\theta_{\text{new}} \leftarrow \theta_{\text{old}} + \Delta \theta_{\text{new}}$$

② RMSprop

$$r_t \leftarrow \rho r_{t-1} + (1-\rho) g_t \otimes g_t \quad (\because g_t = \nabla L(\theta_t))$$

$$\Delta \theta_t \leftarrow -\frac{\epsilon}{\sqrt{s+r_t}} \otimes g_t$$

$$\theta_{t+1} \leftarrow \theta_t + \Delta \theta_t$$

## 19. Adam (adaptive moment estimation)

- Idea: RMSProp + momentum with bias correction

- for each iteration

① compute gradient  $g = \nabla L(\theta)$

② update first moment:  $s \leftarrow \rho_1 s + (1-\rho_1) g$  → momentum

③ update second moment:  $r \leftarrow \rho_2 r + (1-\rho_2) g \otimes g$  → RMSProp

④ bias correction

$$\hat{s} \leftarrow \frac{s}{1-\rho_1^t}$$

$$\hat{r} \leftarrow \frac{r}{1-\rho_2^t}$$

⑤ update parameter

$$\theta \leftarrow \theta - \epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$$

• recommended values in the paper

- learning rate  $\epsilon$ : needs tuning (suggested default: 0.001)
- for momentum  $\rho_1 : 0.9$
- for RMSProp  $\rho_2 : 0.999$
- for stability  $\delta : 10^{-8}$

• Adam: often works better than RMSProp

- recommended as the default algorithm to use
- alternative to Adam worth trying: SGD + Nesterov momentum

### Algorithm 2 Adam optimizer

Require:

- step size  $\epsilon$
- exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1]$
- small constant  $\delta$  used for numerical stabilization
- initial parameters  $\theta$

```

1: initialize 1st and 2nd moment variables  $s = 0$ ,  $r = 0$ 
2: initialize time step  $t = 0$ 
3: while stopping criterion not met do
4:   sample a minibatch  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ 
5:   compute gradient:  $g \leftarrow \frac{1}{m} \nabla_\theta \sum_i L(f(x^{(i)}; \theta), y^{(i)})$ 
6:    $t \leftarrow t + 1$ 
7:   update biased first moment estimate:  $s \leftarrow \rho_1 s + (1-\rho_1) g$ 
8:   update biased second moment estimate:  $r \leftarrow \rho_2 r + (1-\rho_2) g \otimes g$ 
9:   correct bias in first moment:  $\hat{s} \leftarrow \frac{s}{1-\rho_1^t}$ 
10:  correct bias in second moment:  $\hat{r} \leftarrow \frac{r}{1-\rho_2^t}$ 
11:  compute update:  $\Delta \theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$  (operations applied element-wise)
12:  apply update:  $\theta \leftarrow \theta + \Delta \theta$ 
13: end while

```