

# CS 725: OPTIMIZERS

Why do we need more support than what a vanilla SGD learner offers?

Two potential pitfalls:

- (A) Complex (non-convex) loss landscapes that could result in gradients that vary a lot (Soln: Adaptive LR that varies across dimensions)
- (B) Can get stuck in "flat" or "plateau" regions of the loss function (Soln: Momentum)

# ① SGD with momentum

first, define a "velocity" vector

$$V_t \leftarrow \beta V_{t-1} + \eta g_t \quad \text{where} \quad (0 \leq \beta < 1)$$

$$g_t = \nabla_{\omega} L(\omega_{t-1})$$

New weight update rule;  
for SGD w/ momentum

$$\omega_t \leftarrow \omega_{t-1} - V_t$$

$$\begin{aligned}
 \text{Expand } V_t &= \beta V_{t-1} + \eta g_t \\
 &= \eta g_t + \beta(\eta g_{t-1} + \beta V_{t-2}) \\
 &= \eta g_t + \beta \eta g_{t-1} + \beta^2 V_{t-2} \\
 &= \eta g_t + \beta \eta g_{t-1} + \dots + \beta^t V_0
 \end{aligned}$$

What happens to  $V_t$  if  $\beta = 0$ ? Recover Std/vanilla SGD → If  $V_0 = 0$ , then this term disappears

If  $\beta$  is a high value, many past gradients will influence  $w_t$

If  $\beta$  is small  $\approx 0$ , then very few past gradients will influence  $w_t$

## ② RMSProp [for adaptive learning rates]

Consider weights  $w_1, w_2$  where  $w_1 \leftarrow w_1 - \eta \frac{\partial L}{\partial w_1}$   
 $w_2 \leftarrow w_2 - \eta \frac{\partial L}{\partial w_2}$

$$\text{Say } \frac{\partial L}{\partial w_1} \gg \frac{\partial L}{\partial w_2}$$

If  $\eta$  is high, then  $w_1$  can diverge

If  $\eta$  is low, then the update to  $w_2$  will be very small



RMSProp defines a new vector  $S_t$  to adaptively change the learning rate across the weight dimensions

$$S_t \leftarrow \gamma S_{t-1} + (1-\gamma) g_t \odot g_t$$

$(0 \leq \gamma < 1)$  ↗ elementwise product  
or Hadamard product

Expand  $S_t = (1-\gamma) g_t \odot g_t + \gamma S_{t-1}$

$$= (1-\gamma) g_t \odot g_t + \gamma \left( (1-\gamma) g_{t-1} \odot g_{t-1} + \gamma S_{t-2} \right)$$
$$= (1-\gamma) g_t \odot g_t + (1-\gamma) \gamma g_{t-1} \odot g_{t-1} + \dots + \gamma^t S_0$$

What do the coefficients in  $S_t$  add up to?  $\approx 1$  for large  $t$

$S_t$  is an exponentially weighted moving average

New weight update rule:

$$w_t \leftarrow w_{t-1} - \frac{\eta}{\sqrt{S_t + \epsilon}} \odot g_t$$

epsilon is added to avoid  
divide-by-zero errors

③ ADAM optimizer combines ideas from  
SGD w/ momentum and RMS Prop

$$V_t \leftarrow \beta V_{t-1} + (1-\beta) g_t$$

$$S_t \leftarrow \gamma S_{t-1} + (1-\gamma) g_t \odot g_t$$

$$\hat{S}_t \leftarrow \frac{S_t}{1-\gamma^t}$$

$$\hat{V}_t \leftarrow \frac{V_t}{1-\beta^t}$$

Wt update:  $w_t \leftarrow w_{t-1} - \frac{\eta}{\sqrt{\hat{S}_t + \epsilon}} \odot \hat{V}_t$

Needed to correct  
bias in  $S_t$  and  $V_t$  for  
small values of  
 $t$

// BIAS CORRECTION



# CONVOLUTIONAL NEURAL NETWORKS

CNNs are motivated by the visual cortex and aim to satisfy the following desired properties:

- (A) Locality or modeling spatial information
- (B) Translation invariance  $\rightarrow$  model behaviour should be identical regardless of where an object appears in an image
- (C) Parameter-efficient