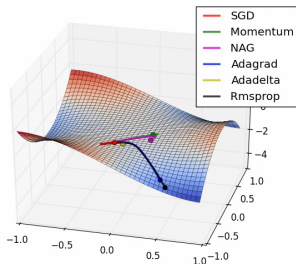# Optimization in Machine Learning

# First order methods: ADAM and friends



**Learning goals**

- Adaptive Step Sizes
- Adagrad
- RMSProp
- ADAM

# ADAPTIVE STEP SIZES

- Step size is probably the most important control param
- Has strong influence on performance
- Natural to use different SS for each input, and to automatically adapt them

# ADAGRAD

- Adagrad adapts SSs by scaling them inversely proportional to square root of the sum of the past squared derivatives
    - Inputs with large partial derivatives get rapid decrease in SS
    - Inputs with small PDs get small decrease in SS
- Goodfellow et al. (2016) say that the accumulation of squared gradients can result in premature decrease in SS

# ADAGRAD

## Algorithm 1 Adagrad

1: **require** Global SS $\alpha$
2: **require** Initial parameter $\boldsymbol{\theta}$
3: **require** Small constant $\beta$, perhaps $10^{-7}$, for numerical stability
4: **Initialize** gradient accumulation variable $\mathbf{r} = \mathbf{0}$
5: **while** stopping criterion not met **do**
6:      Sample a minibatch of $m$ examples from the training set $\{\tilde{x}^{(1)}, \ldots, \tilde{x}^{(m)}\}$
7:      Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L\left(y^{(i)}, f\left(\tilde{\boldsymbol{x}}^{(i)} \mid \boldsymbol{\theta}\right)\right)$
8:      Accumulate squared gradient $\mathbf{r} \leftarrow \mathbf{r} + \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$
9:      Compute update: $\nabla \boldsymbol{\theta} = -\frac{\alpha}{\beta + \sqrt{\mathbf{r}}} \odot \hat{\mathbf{g}}$ (division and square root applied element-wise)
10:      Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \nabla \boldsymbol{\theta}$
11: **end while**

- $\odot$: element-wise product (Hadamard)

# RMSPROP

- Modification of Adagrad
- Resolves Adagrad's radically diminishing SSs.
- Gradient accumulation is replaced by exponentially weighted moving average.
- Theoretically, leads to performance gains in non-convex scenarios.
- Empirically, RMSProp is a very effective optimization algorithm. Particularly, it is employed routinely by DL practitioners.

# RMSPROP

**Algorithm 2** RMSProp

1: **require** Global SS $\alpha$ and decay rate $\rho \in [0, 1)$
2: **require** Initial parameter $\theta$
3: **require** Small constant $\beta$, perhaps $10^{-6}$, for numerical stability
4: Initialize gradient accumulation variable $\mathbf{r} = \mathbf{0}$
5: **while** stopping criterion not met **do**
6:    Sample a minibatch of $m$ examples from the training set $\{\tilde{x}^{(1)}, \ldots, \tilde{x}^{(m)}\}$

7:    Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_\theta \sum_i L\left(y^{(i)}, f\left(\tilde{x}^{(i)} \mid \theta\right)\right)$
8:    Accumulate squared gradient $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$
9:    Compute update: $\nabla \theta = -\frac{\alpha}{\beta + \sqrt{\mathbf{r}}} \odot \hat{\mathbf{g}}$
10:   Apply update: $\theta \leftarrow \theta + \nabla \theta$
11: **end while**

# ADAM

- Adaptive Moment Estimation also has adaptive SSs
- Uses the 1st and 2nd moments of gradients
    - Keeps an exponentially decaying average of past gradients (1st moment)
    - Like RMSProp, stores an exp-decaying avg of past squared gradients (2nd moment)
    - Can be seen as combo of RMSProp + momentum.

# ADAM

## Algorithm 3 Adam

1: **require** Global step size $\alpha$ (suggested default: 0.001)
2: **require** Exponential decay rates for moment estimates, $\rho_1$ and $\rho_2$ in $[0, 1)$ (suggested defaults: 0.9 and 0.999 respectively)
3: **require** Small constant $\beta$ (suggested default $10^{-8}$)
4: **require** Initial parameters $\boldsymbol{\theta}$
5: Initialize time step $t = 0$
6: Initialize 1st and 2nd moment variables $\mathbf{s}^{[0]} = 0, \mathbf{r}^{[0]} = 0$
7: **while** stopping criterion not met **do**
8:      $t \leftarrow t + 1$
9:      Sample a minibatch of $m$ examples from the training set $\{\tilde{x}^{(1)}, \ldots, \tilde{x}^{(m)}\}$
10:      Compute gradient estimate: $\hat{\mathbf{g}}^{[t]} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L\left(y^{(i)}, f\left(\tilde{\boldsymbol{x}}^{(i)} \mid \boldsymbol{\theta}\right)\right)$
11:      Update biased first moment estimate: $\mathbf{s}^{[t]} \leftarrow \rho_1 \mathbf{s}^{[t-1]} + (1 - \rho_1)\hat{\mathbf{g}}^{[t]}$
12:      Update biased second moment estimate: $\mathbf{r}^{[t]} \leftarrow \rho_2 \mathbf{r}^{[t-1]} + (1 - \rho_2)\hat{\mathbf{g}}^{[t]} \odot \hat{\mathbf{g}}^{[t]}$
13:      Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}^{[t]}}{1 - \rho_1^t}$
14:      Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}^{[t]}}{1 - \rho_2^t}$
15:      Compute update: $\nabla \boldsymbol{\theta} = -\alpha \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}}} + \beta}$
16:      Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \nabla \boldsymbol{\theta}$
17: **end while**

# ADAM

- Inits exp-weighted moving averages **s** and **r** as **0** (zero) vectors
- Hence, they are biased towards zero
- This means $\mathbb{E}[\mathbf{s}^{[t]}] \neq \mathbb{E}[\hat{\mathbf{g}}^{[t]}]$ and $\mathbb{E}[\mathbf{r}^{[t]}] \neq \mathbb{E}[\hat{\mathbf{g}}^{[t]} \odot \hat{\mathbf{g}}^{[t]}]$ (where the expectations are calculated over minibatches)
- To see, let's unroll $\mathbf{s}^{[t]}$:

$$\mathbf{s}^{[0]} = 0$$
$$\mathbf{s}^{[1]} = \rho_1 \mathbf{s}^{[0]} + (1 - \rho_1)\hat{\mathbf{g}}^{[1]} = (1 - \rho_1)\hat{\mathbf{g}}^{[1]}$$
$$\mathbf{s}^{[2]} = \rho_1 \mathbf{s}^{[1]} + (1 - \rho_1)\hat{\mathbf{g}}^{[2]} = \rho_1(1 - \rho_1)\hat{\mathbf{g}}^{[1]} + (1 - \rho_1)\hat{\mathbf{g}}^{[2]}$$
$$\mathbf{s}^{[3]} = \rho_1 \mathbf{s}^{[2]} + (1 - \rho_1)\hat{\mathbf{g}}^{[3]} = \rho_1^2(1 - \rho_1)\hat{\mathbf{g}}^{[1]} + \rho_1(1 - \rho_1)\hat{\mathbf{g}}^{[2]} + (1 - \rho_1)\hat{\mathbf{g}}^{[3]}$$

- Therefore, $\mathbf{s}^{[t]} = (1 - \rho_1) \sum_{i=1}^{t} \rho_1^{t-i} \mathbf{g}^{[i]}$.
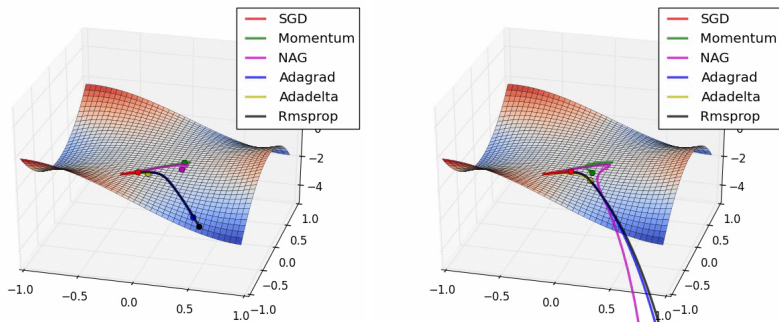- NB: contrib of earlier $\hat{\mathbf{g}}^{[i]}$ to moving average shrinks rapidly

## ADAM

- Now:

$$
\mathbb{E}[\mathbf{s}^{[t]}] = \mathbb{E}[(1 - \rho_1) \sum_{i=1}^{t} \rho_1^{t-i} \hat{\mathbf{g}}^{[i]}]
$$

$$
= \mathbb{E}[\hat{\mathbf{g}}^{[t]}](1 - \rho_1) \sum_{i=1}^{t} \rho_1^{t-i} + \zeta
$$

$$
= \mathbb{E}[\hat{\mathbf{g}}^{[t]}](1 - \rho_1^t) + \zeta
$$

where we approximate $\hat{\mathbf{g}}^{[i]}$ with $\hat{\mathbf{g}}^{[t]}$ which allows us to move it outside the sum. $\zeta$ is the error that results from this approximation.

- Therefore, $\mathbf{s}^{[t]}$ is a biased estimator of $\hat{\mathbf{g}}^{[t]}$ and the effect of the bias vanishes over the time-steps (because $\rho_1^t \to 0$ for $t \to \infty$).

- Ignoring $\zeta$ (as it is small), we correct for the bias by setting $\hat{\mathbf{s}}^{[t]} = \frac{\mathbf{s}^{[t]}}{(1-\rho_1^t)}$.

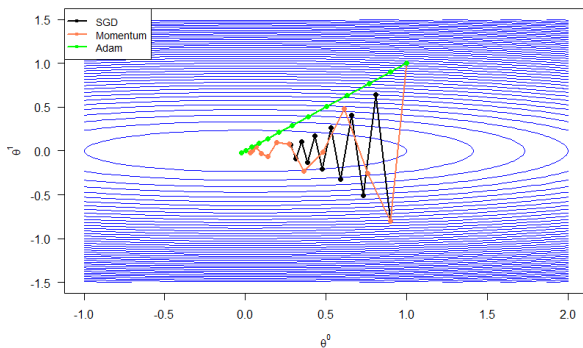- Similarly, we set $\hat{\mathbf{r}}^{[t]} = \frac{\mathbf{r}^{[t]}}{(1-\rho_2^t)}$.

# COMPARISON OF OPTIMIZERS: ANIMATION



Credits: Dettmers (2015) and Radford

Comparison of SGD optimizers near saddle point. Left: After few secs; Right: Later. All methods accelerate compared to vanilla SGD. Best is Rmsprop, then Adagrad.

# COMPARISON ON QUADRATIC FORM



SGD vs. SGD with Momentum vs. ADAM on a quadratic form.