

ADAGRAD / 2

Algorithm AdaGrad

```
1: require Global step size  $\alpha$ 
2: require Initial parameter  $\theta$ 
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 minibatch of  $m$  examples from the training set  $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$ 
7:   Compute gradient estimate:  $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_j L(y^{(j)}, f(\tilde{\mathbf{x}}^{(j)} | \theta))$ 
8:   Accumulate squared gradient  $\mathbf{r} \leftarrow \mathbf{r} + \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$ 
9:   Compute update:  $\nabla \theta \equiv -\frac{\alpha \hat{\mathbf{g}}}{\beta + \sqrt{\mathbf{r}}}$  (operations element-wise)
10:  Apply update:  $\theta \leftarrow \theta + \nabla \theta$ 
11: end while
```

\odot : element-wise product (Hadamard)

\ominus : element-wise product (Hadamard)



RMSPROP / 2

Algorithm RMSProp

```

1: require Global step size  $\alpha$  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 minibatch of  $m$  examples from the training set  $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$ 
7:   Compute gradient estimate:  $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(y^{(i)}, f(\tilde{\mathbf{x}}^{(i)} | \theta))$ 
8:   Accumulate squared gradient  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$ 
9:   Compute update:  $\nabla \theta \equiv -\frac{\beta \alpha}{\sqrt{\mathbf{r}}} \odot \hat{\mathbf{g}}$ 
10:  Apply update:  $\theta \leftarrow \theta + \nabla \theta$ 
11: end while

```



ADAM / 2

Algorithm Adam

- 1: **require** Global step size α (suggested default: 0.001)
- 2: **require** Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$ (suggested defaults: 0.9 and 0.999 respectively)
- 3: **require** Small constant β (suggested default 10^{-8})
- 4: **require** Initial parameters θ
- 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 $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$
- 10: Compute gradient estimate: $\hat{\mathbf{g}}^{[t]} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_j L(y^{(j)}, f(\tilde{\mathbf{x}}^{(j)} | \theta))$
- 11: Update biased first moment estimate: $\hat{\mathbf{s}}^{[t]} \leftarrow \rho_1 \hat{\mathbf{s}}^{[t-1]} + (1 - \rho_1) \hat{\mathbf{g}}^{[t]}$
- 12: Update biased second moment estimate: $\hat{\mathbf{r}}^{[t]} \leftarrow \rho_2 \hat{\mathbf{r}}^{[t-1]} + (1 - \rho_2) \hat{\mathbf{g}}^{[t]} \odot \hat{\mathbf{g}}^{[t]}$
- 13: Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\hat{\mathbf{s}}^{[t]}}{1 - \rho_1^t}$
- 14: Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\hat{\mathbf{r}}^{[t]}}{1 - \rho_2^t}$
- 15: Compute update: $\nabla \theta \equiv -\alpha \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \beta}}$
- 16: Apply update: $\theta \leftarrow \theta + \nabla \theta$
- 17: **end while**

