

Recap: Optimisation

L_2 Regularization

$$\hat{L}_R(\boldsymbol{\theta}) = \hat{L}(\boldsymbol{\theta}) + \frac{\beta}{2} \|\boldsymbol{\theta}\|^2$$

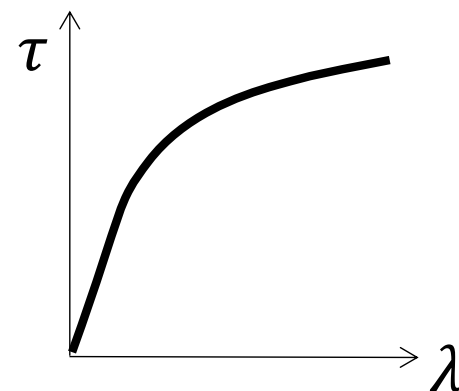
- Effect on gradient descent update

$$\begin{aligned}\boldsymbol{\theta}^{(t+1)} &= \boldsymbol{\theta}^{(t)} - \alpha^{(t)} \nabla_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta}^{(t)}) - \alpha^{(t)} \beta \boldsymbol{\theta}^{(t)} \\ &= (1 - \alpha^{(t)} \beta) \boldsymbol{\theta}^{(t)} - \alpha^{(t)} \nabla_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta}^{(t)})\end{aligned}$$

“weight decay”

- Bayesian interpretation: $\boldsymbol{\theta} \sim \mathcal{N}\left(0, \frac{1}{\sqrt{2\beta}} \mathbf{I}\right)$
- Filter of $\boldsymbol{\theta}^*$** (attenuates along small eigenvectors)

$$\begin{aligned}\boldsymbol{\theta}_R^* &\approx \left(\nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*) + \beta \mathbf{I}\right)^{-1} \nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*) \boldsymbol{\theta}^* \\ &\approx \mathbf{U}(\boldsymbol{\Lambda} + \beta \mathbf{I})^{-1} \boldsymbol{\Lambda} \mathbf{U}^T \boldsymbol{\theta}^* \\ &= \mathbf{U} \boldsymbol{\tau}(\boldsymbol{\Lambda}) \mathbf{U}^T \boldsymbol{\theta}^* \\ &= \boldsymbol{\tau}\left(\nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*)\right) \boldsymbol{\theta}^*\end{aligned}$$



Noisy input

- Add i.i.d. noise $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \beta \mathbf{I})$ to input
- Loss becomes

$$\begin{aligned} L(\mathbf{w}) &= \mathbb{E}(\mathbf{w}^T(\mathbf{x} + \boldsymbol{\epsilon}) - y)^2 \\ &= \mathbb{E}(\mathbf{w}^T \mathbf{x} - y)^2 + 2\mathbb{E} \mathbf{w}^T \boldsymbol{\epsilon} (\mathbf{w}^T \mathbf{x} - y) + \mathbb{E}(\mathbf{w}^T \boldsymbol{\epsilon})^2 \\ &= \mathbb{E}(\mathbf{w}^T \mathbf{x} - y)^2 + \mathbf{w}^T \mathbb{E}(\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T) \mathbf{w} \\ &= \mathbb{E}(\mathbf{w}^T \mathbf{x} - y)^2 + \beta \|\mathbf{w}\|^2 \end{aligned}$$

- Equivalent to weight decay!

Descent method: general recipe

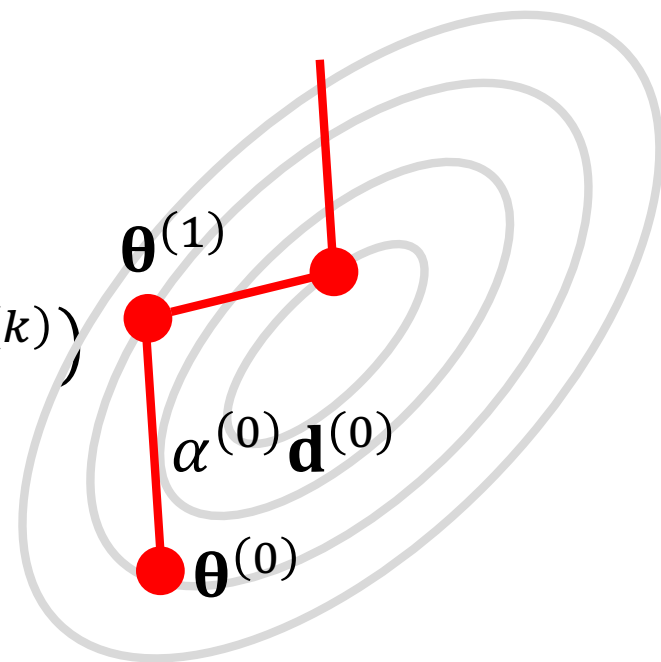
Initialization: start with some $\boldsymbol{\theta}^{(0)}$

For $k = 0, \dots$ **until convergence**

Choose **descent direction** $\mathbf{d}^{(k)} = -\nabla \hat{L}(\boldsymbol{\theta}^{(k)})$

Choose **step size** $\alpha^{(k)}$

Update $\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + \alpha^{(k)} \mathbf{d}^{(k)}$



Gradient descent convergence rate

Strong convexity $\nabla^2 \hat{L}(\boldsymbol{\theta}) \succcurlyeq m\mathbf{I} \quad m > 0$

Lipschitz gradient $\nabla^2 \hat{L}(\boldsymbol{\theta}) \preccurlyeq M\mathbf{I}$

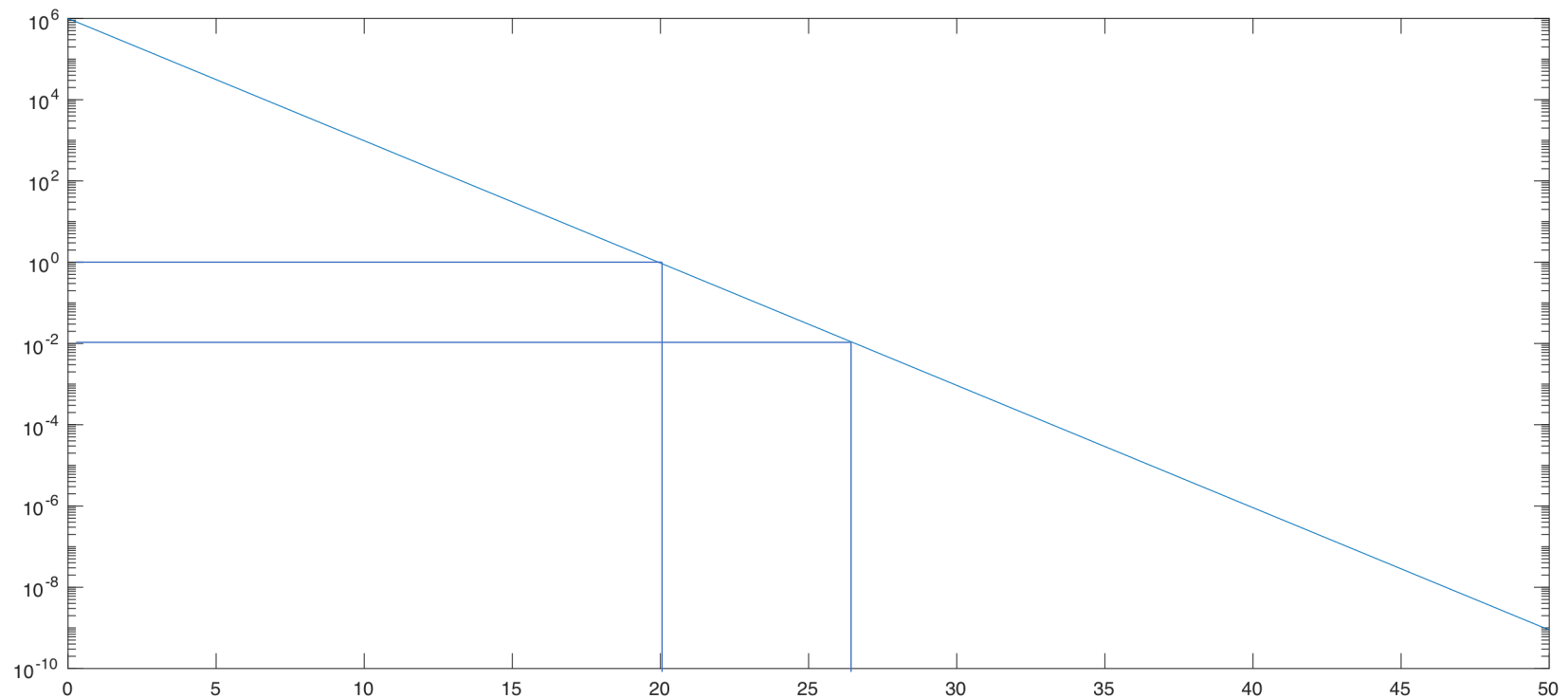
Constant step size $\alpha \leq \frac{2}{m + M}$

$$\hat{L}(\boldsymbol{\theta}^{(k)}) - \hat{L}(\boldsymbol{\theta}^*) \leq c^k \cdot \frac{M}{2} \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^{(k)}\|$$

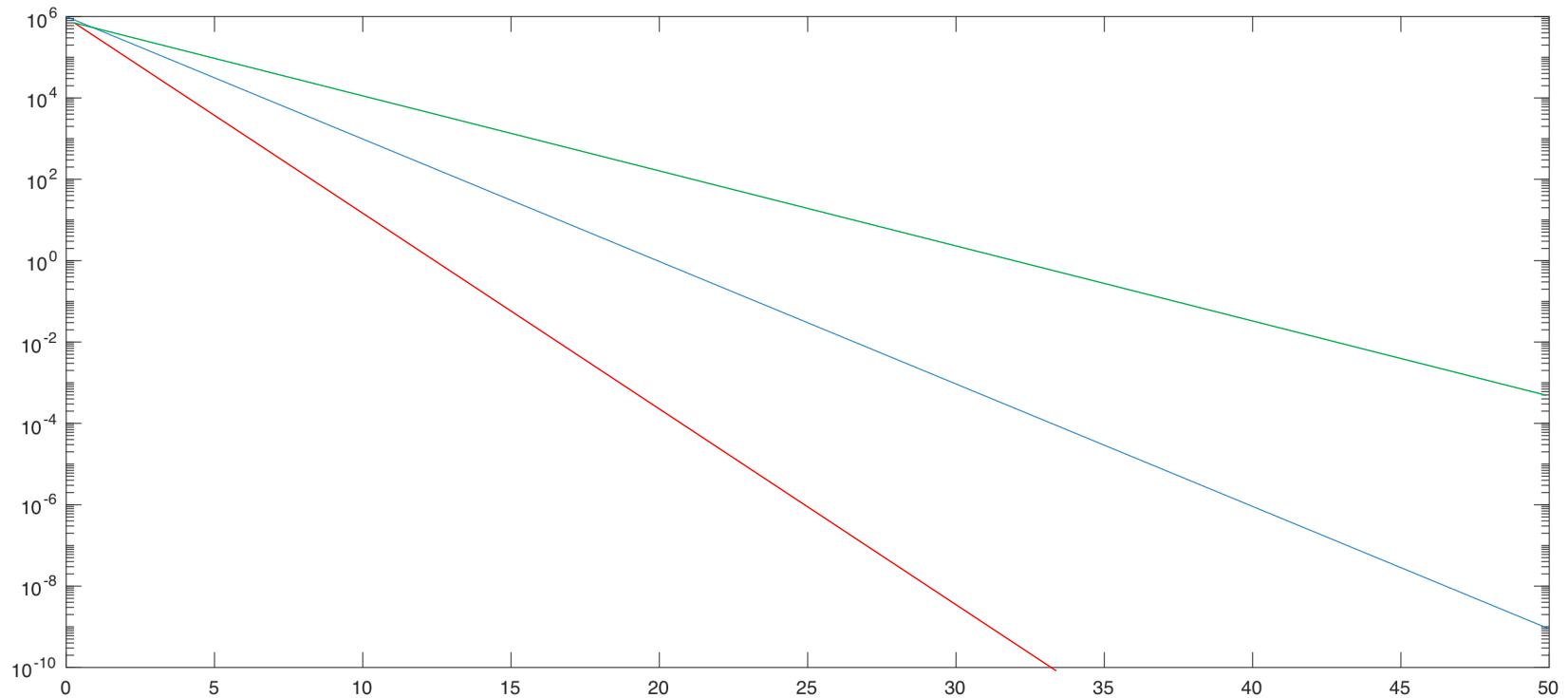
“Linear” convergence

To get $\hat{L}(\boldsymbol{\theta}^{(k)}) - \hat{L}(\boldsymbol{\theta}^*) \leq \epsilon$ one needs $\mathcal{O}\left(\log \frac{1}{\epsilon}\right)$ iterations

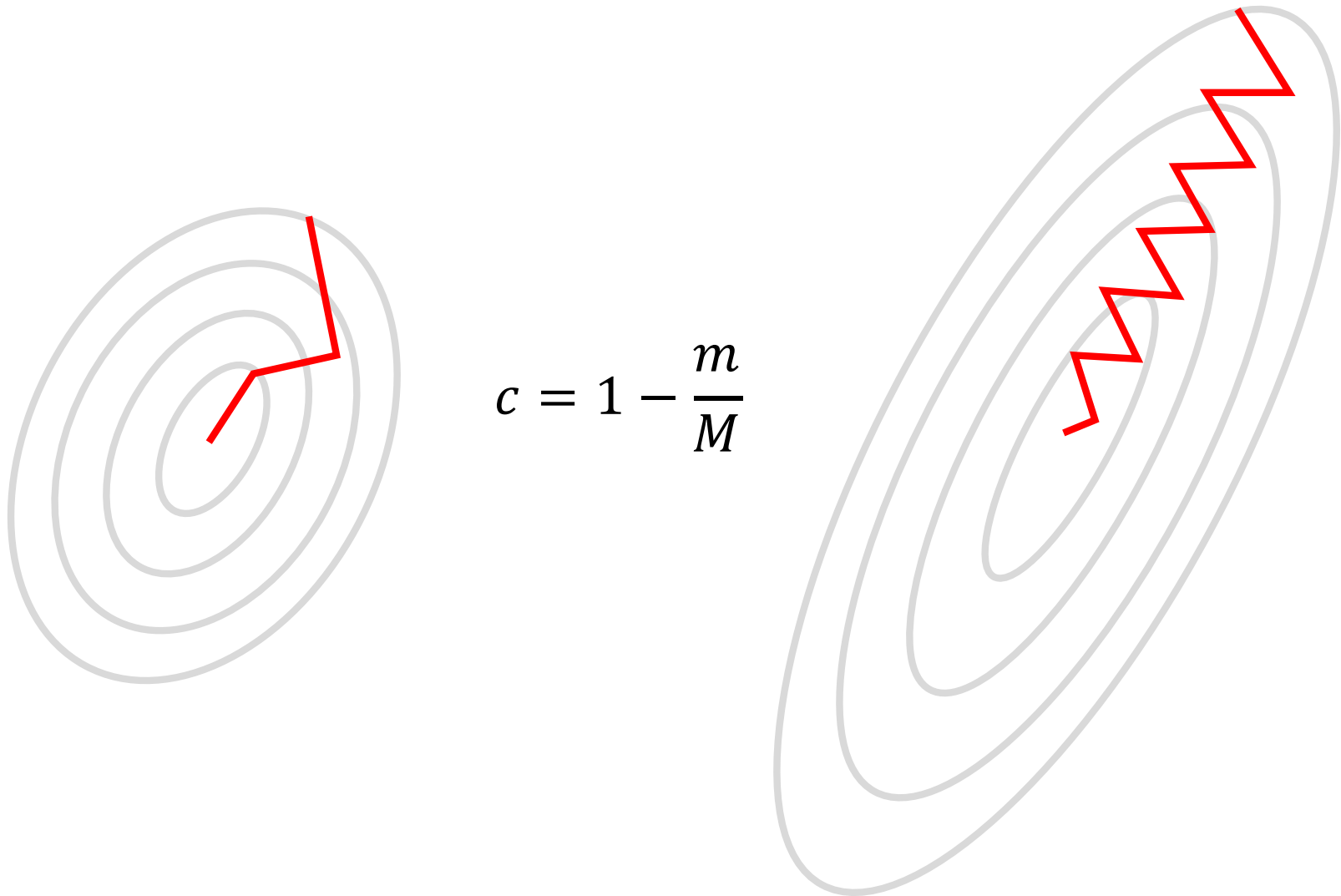
Gradient descent convergence rate



Gradient descent convergence rate



Gradient descent convergence rate



Momentum

- Key idea: add back fraction of past step

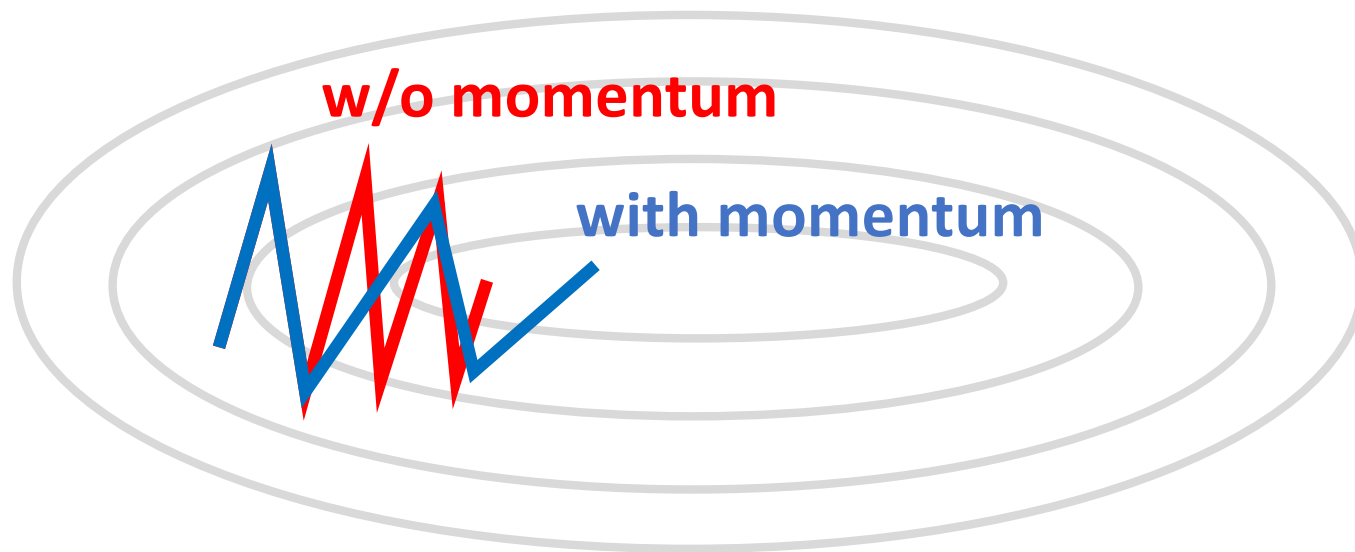
$$\Delta \boldsymbol{\theta}^{(k)} = -\alpha \nabla \hat{L}(\boldsymbol{\theta}^{(k)}) + \beta \Delta \boldsymbol{\theta}^{(k-1)} \quad 0 < \beta < 1$$
$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} + \Delta \boldsymbol{\theta}^{(k)}$$

- Exponentially decaying memory

$$\Delta \boldsymbol{\theta}^{(k)} = -\alpha \left[\nabla \hat{L}(\boldsymbol{\theta}^{(k)}) + \beta \nabla \hat{L}(\boldsymbol{\theta}^{(k-1)}) + \dots \beta^n \nabla \hat{L}(\boldsymbol{\theta}^{(k-n)}) \right] \\ + \beta^{n+1} \Delta \boldsymbol{\theta}^{(k-n-1)}$$

- Dampen oscillations of gradient descent

Momentum



AdaGrad

- Key idea: use different update per coordinate

$$\theta_i^{(k+1)} = \theta_i^{(k)} - \frac{\alpha}{\sqrt{g_{ii}^{(k)} + \epsilon}} \frac{\partial}{\partial \theta_i} \hat{L}(\boldsymbol{\theta}^{(k)})$$

$$\mathbf{G}^{(k)} = \sum_{l=1}^k \nabla \hat{L}(\boldsymbol{\theta}^{(l)}) \nabla \hat{L}(\boldsymbol{\theta}^{(l)})^T$$

Covariance up to
scaling factor

AdaGrad

- Key idea: use different update per coordinate

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha (\text{Diag} \mathbf{G}^{(k)} + \epsilon \mathbf{I})^{-1/2} \nabla \hat{L}(\boldsymbol{\theta}^{(k)})$$

Matrix without
off-diag elements

$$\mathbf{G}^{(k)} = \sum_{l=1}^k \nabla \hat{L}(\boldsymbol{\theta}^{(l)}) \nabla \hat{L}(\boldsymbol{\theta}^{(l)})^T$$

- Eliminates the need to manually tune learning rate
- Accumulated sum in denominator makes learning rate too small

AdaDelta or RMSProp

- Modification of AdaGrad using decaying average of gradients

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \frac{\alpha}{\sqrt{\mathbf{g}^{(k)} + \epsilon}} \nabla \hat{L}(\boldsymbol{\theta}^{(k)})$$

$$\mathbf{g}^{(k)} = \gamma \mathbf{g}^{(k-1)} + (1 - \gamma) \nabla \hat{L}^2(\boldsymbol{\theta}^{(k)})$$

* All operations performed coordinate-wise
Zeiler 2012; Hinton (unpublished)

Adam (Adaptive Moment Estimation)

- Key idea: keep decaying average of gradients and their squares

$$\mathbf{m}^{(k)} = \beta_1 \mathbf{m}^{(k-1)} + (1 - \beta_1) \nabla \hat{L}(\boldsymbol{\theta}^{(k)})$$

$$\mathbf{v}^{(k)} = \beta_2 \mathbf{v}^{(k-1)} + (1 - \beta_2) \nabla \hat{L}^2(\boldsymbol{\theta}^{(k)})$$

$$\hat{\mathbf{m}}^{(k)} = \frac{\mathbf{m}^{(k)}}{1 - \beta_1^k} \qquad \hat{\mathbf{v}}^{(k)} = \frac{\mathbf{v}^{(k)}}{1 - \beta_2^k}$$

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha \frac{\hat{\mathbf{m}}^{(k)}}{\sqrt{\hat{\mathbf{v}}^{(k)} + \epsilon}}$$