

## Exponential Moving Averages

$$V_0 = 0$$

$$V_1 = \beta * V_0 + (1-\beta) * \Theta_1 \text{ [Current observation at time 1] } [\beta=0.9]$$

$$V_2 = \beta * V_1 + (1-\beta) * \Theta_2 \text{ [Current observation at time 2]}$$

$$V_t = \beta * V_{t-1} + (1-\beta) * \Theta_t \text{ [Current observation at time t]}$$

Generalized form

$$V = \beta * V + (1-\beta) * \Theta$$

## Gradient Descent (mini-batch) with Momentum (SGD - Minibatch)

For each batch:

Compute dw, db

$$V_{dw} = \beta * V_{dw} + (1-\beta) * dw$$

$$V_{db} = \beta * V_{db} + (1-\beta) * db$$

$$w = w - \alpha * V_{dw}$$

$$b = b - \alpha * V_{db}$$

## Root Mean Square Propagation (RMSProp)

$$S_{dw} = \beta^2 * S_{dw} + (1-\beta^2) * dw^2$$

$$S_{db} = \beta^2 * S_{db} + (1-\beta^2) * db^2$$

$$w = w - \alpha * dw / \sqrt{S_{dw} + \epsilon}$$

$$b = b - \alpha * db / \sqrt{S_{db} + \epsilon}$$

Epsilon = small value to prevent division by 0. Normally  $10^{-8}$

## Adaptive Moment with Estimation (Adam)

For each batch:

$$V_{dw} = \beta_1 * V_{dw} + (1-\beta_1) * dw$$

$$V_{dw} = V_{dw} / (1-\beta_1^t)$$

$$V_{db} = \beta_1 * V_{db} + (1-\beta_1) * db$$

$$V_{db} = V_{db} / (1-\beta_1^t)$$

$$S_{dw} = \beta_2 * S_{dw} + (1-\beta_2) * dw^2$$

$$S_{dw} = S_{dw} / (1-\beta_2^t)$$

$$S_{db} = \beta_2 * S_{db} + (1-\beta_2) * db^2$$

$$S_{db} = S_{db} / (1-\beta_2^t)$$

$$w = w - \alpha * V_{dw} / [\sqrt{S_{dw}} + \epsilon]$$

$$b = b - \alpha * V_{db} / [\sqrt{S_{db}} + \epsilon]$$

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

$$\epsilon = 10^{-8}$$