**Week 2-2 Adam Algo, GC with Momentum (Exponentially Weight MA) + RMS prop**

# adam paper: https://arxiv.org/pdf/1412.6980.pdf

# Exponentially weighted moving averages



Temperature in London

$\theta_1 = 40°F$  4°C

$\theta_2 = 49°F$  9°C

$\theta_3 = 45°F$

$\vdots$

$\theta_{180} = 60°F$  15°C

$\theta_{181} = 56°F$

$\vdots$

$V_0 = 0$

$V_1 = 0.9 V_0 + 0.1 \theta_1$

$V_2 = 0.9 V_1 + 0.1 \theta_2$

$V_3 = 0.9 V_2 + 0.1 \theta_3$

$\vdots$

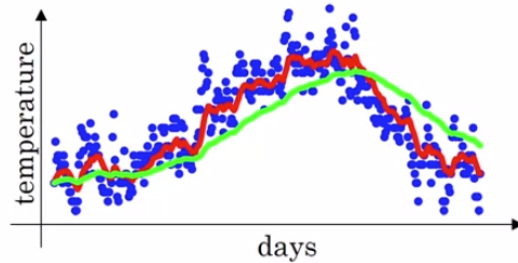$V_t = 0.9 V_{t-1} + 0.1 \theta_t$

# Exponentially weighted averages

$$V_t = \beta V_{t-1} + (1-\beta)\theta_t$$

$\beta = 0.9$ : ≈ 10 days' temperature

$\beta = 0.98$ : ≈ 50 days

$V_t$ is approximately average over ≈ $\frac{1}{1-\beta}$ days' temperature.

$$\frac{1}{1-0.98} = 50$$


temperature vs days

# Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1-\beta)\theta_t$$
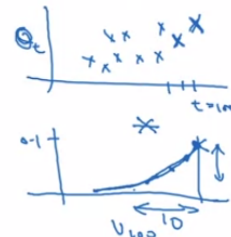
$$v_{100} = 0.9v_{99} + 0.1\theta_{100}$$
$$v_{99} = 0.9v_{98} + 0.1\theta_{99}$$
$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$

...

$\to V_{100} = 0.1\,\theta_{100} + 0.9 \times (0.1\,\theta_{99} + 0.9 \times v_{98})$

$0.1\theta_{98} + 0.9 v_{97}$

$= 0.1\,\theta_{100} + 0.1\times 0.9\cdot\theta_{99} + 0.1\,(0.9)^2\theta_{98} + 0.1\,(0.9)^3\theta_{97} + 0.1\,(0.9)^4\theta_{96}$
$+ \cdots$

$0.9^{10} \approx 0.35 \approx \frac{1}{e}$

$(1-\varepsilon)^{1/\varepsilon} = \frac{1}{e}$ , $\varepsilon = 0.9$

$0.98^{?}$

$\to 0.98^{50} \approx \frac{1}{e}$

(all these weights approximately add up to 1)
after 1/epi -> 1/e weight, rather small

## Implementing exponentially weighted averages

$$v_0 = 0$$
$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$
$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$
$$v_3 = \beta v_2 + (1 - \beta) \theta_3$$
...

$$v_\theta := 0$$
$$v_\theta := \beta v + (1-\beta)\theta_1$$
$$v_\theta := \beta v + (1-\beta)\theta_2$$
$$\vdots$$

$\rightarrow v_\theta = 0$

Repeat {

 Get next $\theta_t$

 $v_\theta := \beta v_\theta + (1-\beta)\theta_t$ $\leftarrow$
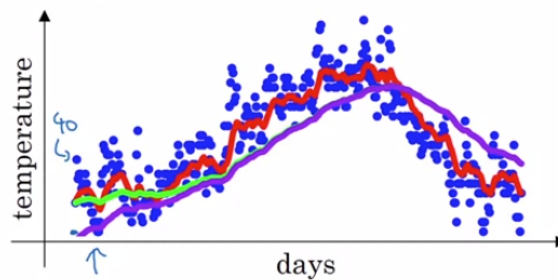
}

Andrew Ng

(memory efficient)

---

# Bias correction

## Bias correction



$\rightarrow v_t = \beta v_{t-1} + (1-\beta)\theta_t$

$v_0 = 0$

$v_1 = \cancel{0.98 v_0} + 0.02\,\theta_1$

$v_2 = 0.98\, v_1 + 0.02\,\theta_2$

$\quad = 0.98 \times 0.02 \times \theta_1 + 0.02\,\theta_2$

$\quad = 0.0196\,\theta_1 + 0.02\,\theta_2$

$\beta = 0.98$

Andrew Ng

if we initialize with 0, it gives the purple curve

## Bias correction



$\beta = 0.98$

$$\rightarrow v_t = \beta v_{t-1} + (1-\beta)\theta_t$$

$V_0 = 0$

$V_1 = \cancel{0.98 V_0} + 0.02\,\theta_1$
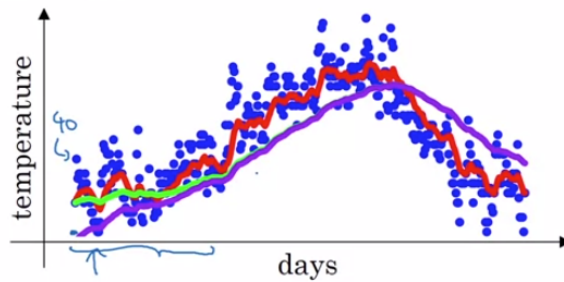
$V_2 = 0.98\,V_1 + 0.02\,\theta_2$

$\quad = 0.98 \times 0.02 \times \theta_1 + 0.02\,\theta_2$

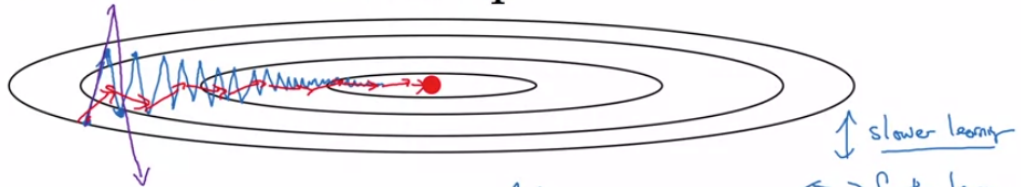$\quad = 0.0196\,\theta_1 + 0.02\,\theta_2$

$\dfrac{V_t}{1-\beta^t}$

$t=2:\quad 1-\beta^t = 1-(0.98)^2 = 0.0396$

$\dfrac{V_2}{0.0396} = \dfrac{0.0196\,\theta_1 + 0.02\,\theta_2}{0.0396}$

Andrew Ng

# average of these derivatives

## Gradient descent example



↑ slower learning

⟷ faster learning

Momentum:

On iteration $t$:

Compute $dW, db$ on current mini-batch.

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

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

"$V_\theta = \beta V_{\theta\tau}(1-\beta)\theta_t$"

$W := W - \alpha V_{dW}$ , $b := b - \alpha V_{db}$

## Gradient descent example



↑ slower learning

⟷ faster learning

Momentum:

On iteration $t$:

Compute $dW, db$ on current mini-batch.

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

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

friction ⎯⎯ ↕ velocity ⎯⎯ ↑ acceleration

"$V_\theta = \beta V_{\theta\tau}(1-\beta)\theta_t$"

$W := W - \alpha V_{dW}$ , $b := b - \alpha V_{db}$

(add friction so that its horizontal speed along the bowl gets smaller)

## Implementation details

$v_{dw} = 0$

On iteration $t$:

    Compute $dW, db$ on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta)dW$$
$$v_{db} = \beta v_{db} + (1 - \beta)db$$
$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

$\frac{v_{dw}}{1-\beta^t}$

Hyperparameters: $\alpha, \beta$      $\beta = 0.9$

       ↑ ↑           average over last ≈ 10 gradients
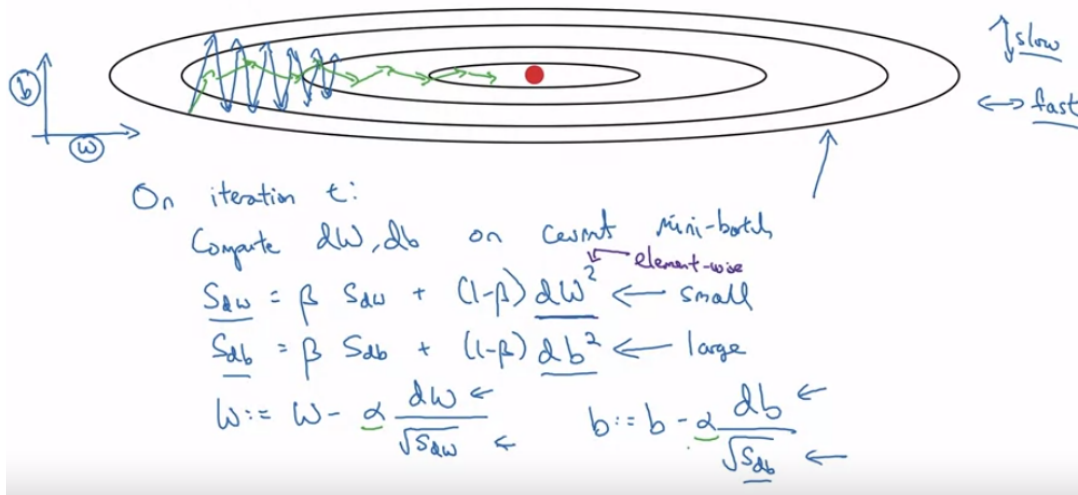
In practice, people don't usually do this (bias correction) because after just ten iterations, your moving average will have warmed up and is no longer a bias estimate.

We sometimes omit the (1-beta) term, but it is not so intuitive and we may need to adjust the learning rate correspondingly

RMSprop

# RMSprop



On iteration $t$:

Compute $dW, db$ on current mini-batch

$S_{dW} = \beta \, S_{dW} + (1-\beta) \, dW^2 \leftarrow$ small ← element-wise

$S_{db} = \beta \, S_{db} + (1-\beta) \, db^2 \leftarrow$ large

$W := W - \alpha \dfrac{dW}{\sqrt{S_{dW}}}$

$b := b - \alpha \dfrac{db}{\sqrt{S_{db}}}$

larger db gets more scale division, also enables larger learning rate

We also add a epsilon for numerical stability (division by zero)

# RMSprop



On iteration $t$:

Compute $dW, db$ on current mini-batch

$S_{dW} = \beta_2 \, S_{dW} + (1-\beta_2) \, dW^2 \leftarrow$ small ← element-wise

$S_{db} = \beta_2 \, S_{db} + (1-\beta_2) \, db^2 \leftarrow$ large

$W := W - \alpha \dfrac{dW}{\sqrt{S_{dW}} + \varepsilon}$

$b := b - \alpha \dfrac{db}{\sqrt{S_{db}} + \varepsilon}$

$\varepsilon = 10^{-8}$

Andrew Ng

# Adam

## Adam optimization algorithm

$V_{dw} = 0, S_{dw} = 0. \quad V_{db} = 0, S_{db} = 0$

On iteration $t$:

Compute $dw, db$ using current mini-batch

$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dW \quad, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db \quad \leftarrow \text{"momentum"} \; \beta_1$

$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2 \quad, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db \quad \leftarrow \text{"RMSprop"} \; \beta_2$

$V_{dw}^{corrected} = V_{dw}/(1-\beta_1^t) \quad, \quad V_{db}^{corrected} = V_{db}/(1-\beta_1^t)$

$S_{dw}^{corrected} = S_{dw}/(1-\beta_2^t) \quad, \quad S_{db}^{corrected} = S_{db}/(1-\beta_2^t)$

$W := W - \alpha \dfrac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected}} + \varepsilon} \qquad b := b - \alpha \dfrac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected}} + \varepsilon}$

(there is db^2)

## Hyperparameters choice:

$\rightarrow \alpha$ : needs to be tune

$\rightarrow \beta_1$ : 0.9 $\qquad \rightarrow (dw)$

$\rightarrow \beta_2$ : 0.999 $\qquad \rightarrow (dw^2)$

$\rightarrow \varepsilon$ : $10^{-8}$

Adam : Adaptive moment estimation

Adaptive moment estimation (first order momentum + second order momentum)