



Deep Learning

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir



<https://www.aparat.com/mehran.safayani>



https://github.com/safayani/deep_learning_course



Optimization Algorithms

Understanding exponentially weighted averages

Exponential moving average

$$v_0 = 0$$

$$v_1 = 0.9v_0 + 0.1\theta_1$$

$$v_2 = 0.9v_1 + 0.1\theta_2$$

$$v_3 = 0.9v_2 + 0.1\theta_3$$

.

.

.

$$v_t = 0.9v_{t-1} + 0.1\theta_t$$

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

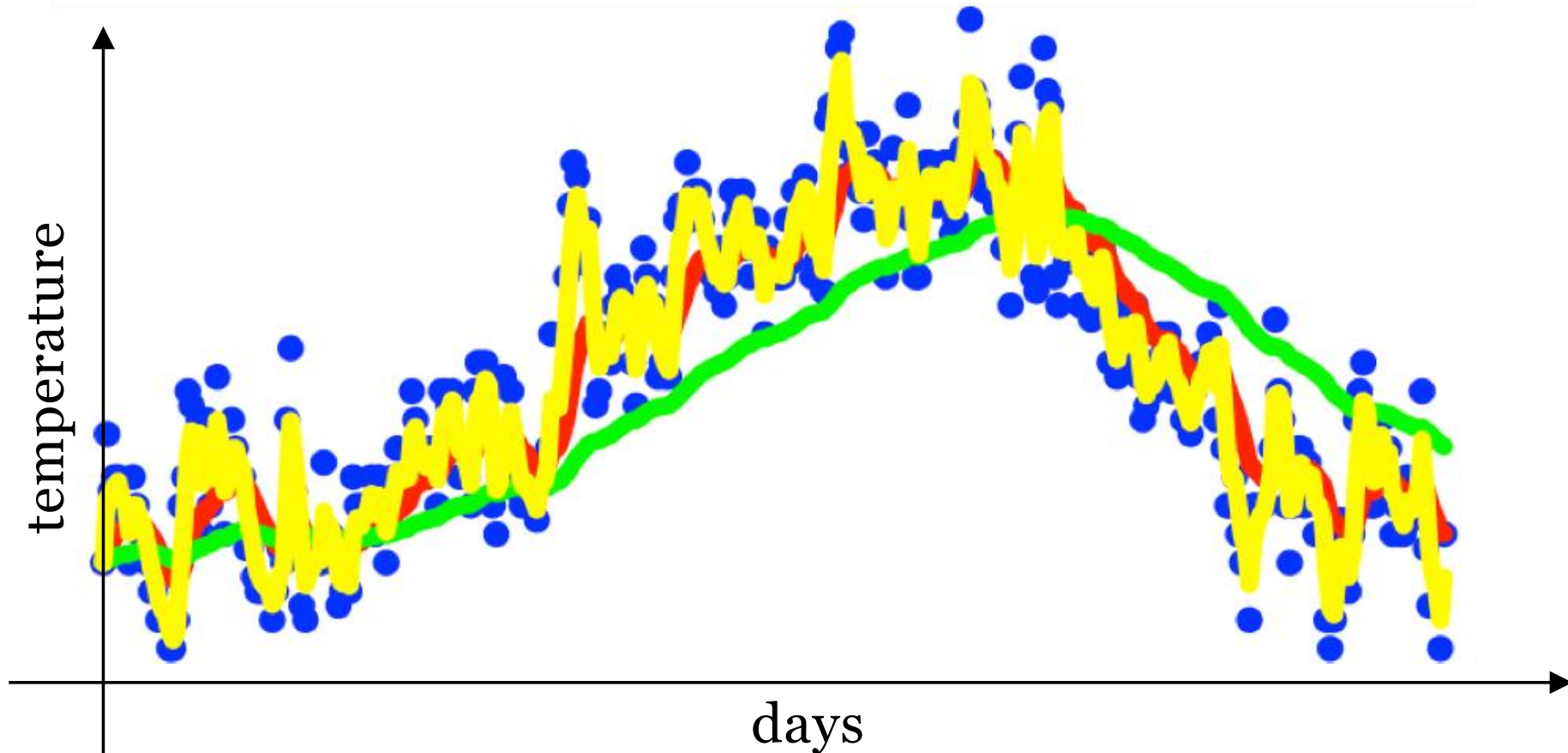
$$\beta = 0.9 \approx 10 \text{ days} \quad \frac{1}{1-\beta} \text{ days}$$

$$\beta = 0.98 \approx 50 \text{ days} \quad \frac{1}{1-0.98} = 50$$

$$\beta = 0.5 \approx 2 \text{ days}$$

Exponentially weighted averages

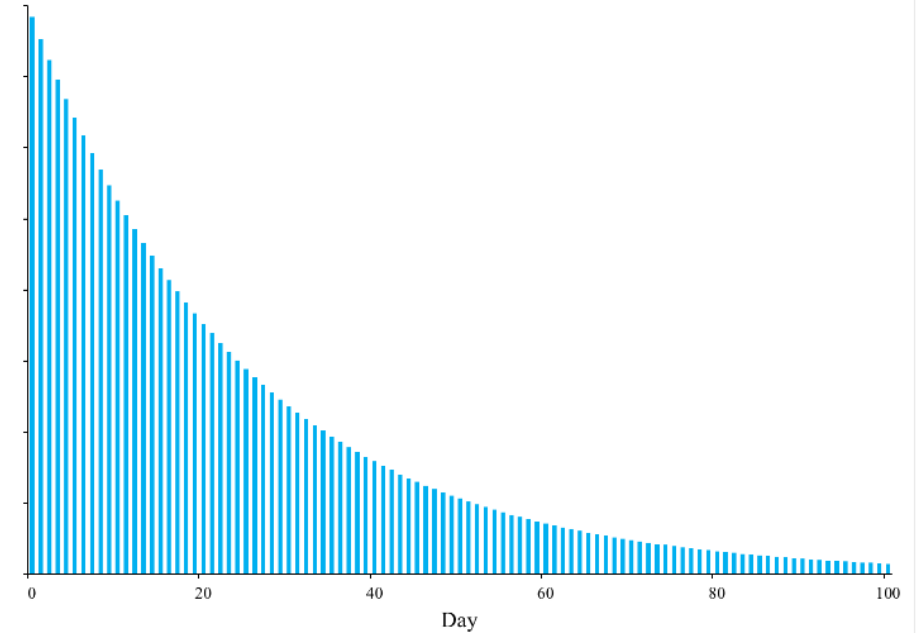
- $v_t = \beta v_{t-1} + (1 - \beta)\theta_t$



Exponentially weighted averages

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

.....



- $V_{100} = 0.1\theta_{100} + 0.9(0.9V_{98} + 0.1\theta_{99}) = \mathbf{0.1\theta_{100} + 0.1 \times 0.9\theta_{99} + 0.1 \times (0.9)^2\theta_{98} + 0.1 \times (0.9)^3\theta_{97} + 0.1 \times (0.9)^4\theta_{96}}$

$$\text{\#days} = \frac{1}{1 - \beta}$$

$$(0.9)^{10} \cong 0.35 = \frac{1}{e}$$

$$\varepsilon = 1 - \beta$$

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

$$((1 - 0.1)^{1/0.1}) = \frac{1}{e}$$

$$\beta = 0.98$$

$$((1 - 0.02)^{1/0.02}) \cong e \cong 0.35$$

Implementing exponentially weighted averages

$$V_{\theta}=0$$

$$V_{\theta}= V_{\theta} + (1 - \beta)\theta_1$$

$$V_{\theta}= V_{\theta} + (1 - \beta)\theta_2$$

$$V_{\theta}=0$$

Repeat {

 Get Next θ_t

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

}

Bias correction in exponentially weighted average

- **Bias correction:**

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

$$V_0 = 0$$

$$V_1 = 0.98V_0 + 0.02\theta_1$$

$$V_2 = 0.98V_1 + 0.02\theta_2$$

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

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



$$V_t = \frac{V_t}{1 - \beta^t}$$

$$t=2$$

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

$$\tilde{V}_2 = \frac{V_2}{0.0396} = \frac{0.0196\theta_1 + 0.02\theta_2}{0.0396}$$

$$V_t = \frac{V_t}{1 - \beta^t}$$



بزرگ شود t \rightarrow $\beta^t = 0$

Optimization Algorithms

Gradient descent with momentum

- Momentum:

On iteration t :

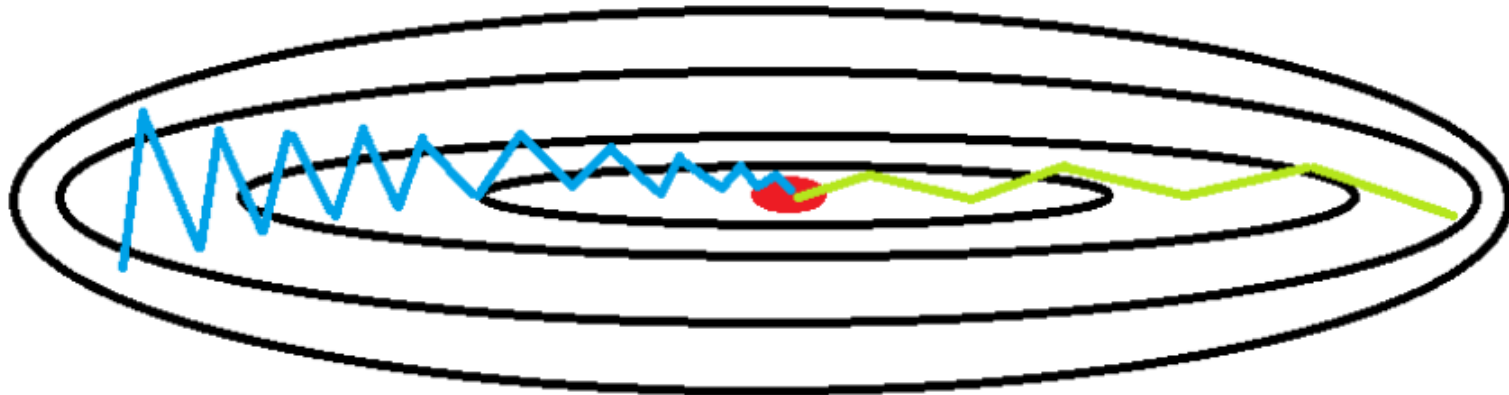
Compute $d w$, $d b$ on the current mini-batch

$$v_{d w} = \beta v_{d w} + (1 - \beta) d w$$

$$v_{d b} = \beta v_{d b} + (1 - \beta) d b$$

$$w = w - \alpha v_{d w}, b = b - \alpha v_{d b}$$

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



Gradient descent with momentum

Implementation details

$$v_{dw} = 0, v_{db} = 0$$

On Iteration t:

$$v_{dw} = \beta v_{dw} + (1 - \beta)dw \equiv v_{dw} = \hat{\beta} v_{dw} + dw$$

$$v_{db} = \hat{\beta} v_{db} + db$$

$$w = w - \alpha v_{dw}$$

$$b = b - \alpha v_{db}$$

$$\alpha [\beta v_{dw} + (1 - \beta)dw] = \underbrace{(1 - \beta)\alpha}_{\boxed{\alpha}} \underbrace{\left[\frac{\beta}{1 - \beta} v_{dw} + dw\right]}_{\boxed{\hat{\beta}}} = \hat{\alpha} [\hat{\beta} v_{dw} + dw]$$

Optimization Algorithms

RMSprop(Root Mean Square Propagation)

- RMS Prop:

On iteration t:

Compute dw , db on current mini-batch

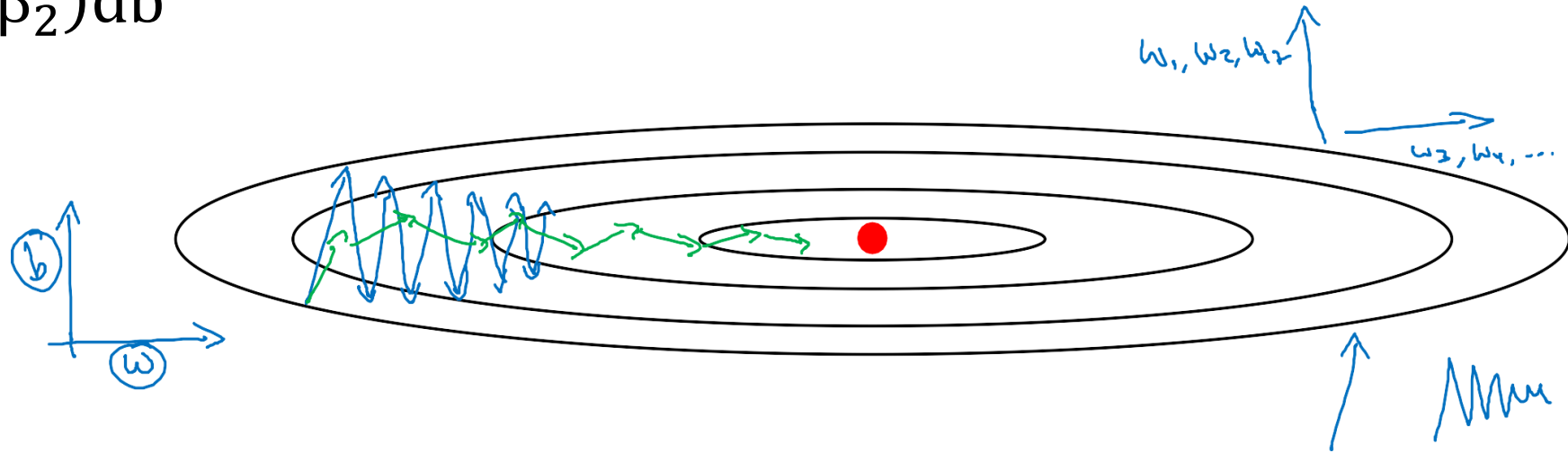
$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2 \longrightarrow \text{Element wise}$$

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

$$W = W - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}$$

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

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



Optimization Algorithms

Adam: Adaptive momentum estimation

- Adam:

$$v_{dw} = 0, S_{db} = 0, v_{db} = 0, S_{db} = 0$$

On iteration t :

Compute dw , db using current mini-batch

momentum

$$v_{dw} = \beta_1 v_{dw} + (1 - \beta_1)dw, v_{db} = \beta_1 v_{db} + (1 - \beta_1)db$$

RMS Prop

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

$$v_{dw}^{corrected} = \frac{v_{dw}}{1 - \beta_1^t}, v_{db}^{corrected} = \frac{v_{db}}{1 - \beta_1^t}$$

$$S_{dw}^{corrected} = \frac{S_{dw}}{1 - \beta_2^t}, S_{db}^{corrected} = \frac{S_{db}}{1 - \beta_2^t}$$

$$w = w - \alpha \frac{v_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}}, b = b - \alpha \frac{v_{db}^{corrected}}{\sqrt{S_{db}^{corrected} + \epsilon}}$$

α = needs to be tuned

$\beta_1 = 0.9$

$\beta_2 = 0.999$

$\epsilon = 10^{-8}$

Optimization Algorithms

Learning rate decay

- $$\alpha = \frac{1}{1 + (\text{decay-rate}) * (\text{epoch-num})} \alpha_0$$

Epoch	α
1	0.01
2	0.006
3	0.005
4	0.004
...	...

$$\alpha_0 = 0.02$$

$$\text{decay-rate} = 1$$

- $$\alpha = (0.95)^{\# \text{epoch-num}} \times \alpha_0$$

- $$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \times \alpha_0$$



$$\alpha = \frac{1}{\sqrt{t}} \times \alpha_0$$

#mini-batch