



deeplearning.ai

Optimization Algorithms

Mini-batch gradient descent

entire training set
for each iteration

a mini-batch
for each iteration

Batch vs. mini-batch gradient descent

Vectorization allows you to efficiently compute on m examples.

$$X = \underbrace{[x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(500)}]}_{\substack{(n_x, m) \\ X^{\{1\}} \ (n_x, 1000)}} \left[\begin{array}{c} x^{(1001)} \\ \vdots \\ x^{(2000)} \end{array} \right] \ \dots \ \left[\begin{array}{c} \dots \\ \dots \\ x^{(m)} \end{array} \right] \ X^{\{2\}} \ (n_x, 1000) \ \dots \ X^{\{5,000\}} \ (n_x, 1000)$$

$$Y = \underbrace{[y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)}]}_{\substack{(1, m) \\ Y^{\{1\}} \ (1, 1000)}} \left[\begin{array}{c} y^{(1001)} \\ \vdots \\ y^{(2000)} \end{array} \right] \ \dots \ \left[\begin{array}{c} \dots \\ \dots \\ y^{(m)} \end{array} \right] \ Y^{\{2\}} \ (1, 1000) \ \dots \ Y^{\{5,000\}} \ (1, 1000)$$

What if $m = 5,000,000$?

5,000 mini-batches of 1,000 each

Mini-batch t : $\underline{X^{\{t\}}, Y^{\{t\}}}$

$x^{(i)}$ i -th training example
 $z^{[l]}$ l -th layer
 $X^{\{t\}}, Y^{\{t\}}$ t -th minibatch

$$\begin{cases} X = (n_x, m) \Rightarrow X^{\{t\}} = (n_x, \frac{m}{k}) \\ Y = (1, m) \Rightarrow Y^{\{t\}} = (1, \frac{m}{k}) \end{cases} \quad k = \# \text{mini batches}$$

Algorithm

Mini-batch gradient descent

Repeat {
for $t = 1, \dots, 5000$ {

→ Forward prop on $X^{[t]}$.

$$\begin{aligned} Z^{[t]} &= W^{[t]} X^{[t]} + b^{[t]} \\ A^{[t]} &= g^{[t]}(Z^{[t]}) \\ &\vdots \\ A^{[t]} &= g^{[t]}(Z^{[t]}) \end{aligned}$$

} Vectorized implementation
(1000 examples)

→ Compute cost $J^{[t]} = \frac{1}{1000} \sum_{i=1}^l \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|W^{[l]}\|_F^2$.

Backprop to compute gradients w.r.t $J^{[t]}$ (using $(X^{[t]}, Y^{[t]})$)

$$W^{[t]} := W^{[t]} - \alpha \delta W^{[t]}, \quad b^{[t]} := b^{[t]} - \alpha \delta b^{[t]}$$

} }

One step of gradient descent
using $\frac{X^{[t]}, Y^{[t]}}{(as if n=1000)}$

Instead of $\underline{X}, \underline{Y}$ use X, Y

"1 epoch"
pass through train set.
a single



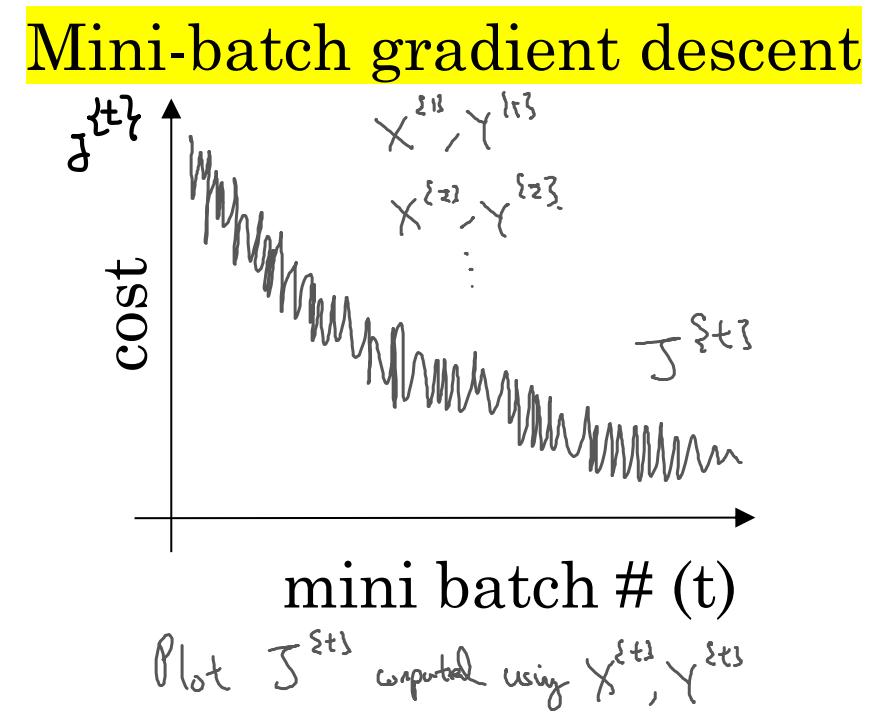
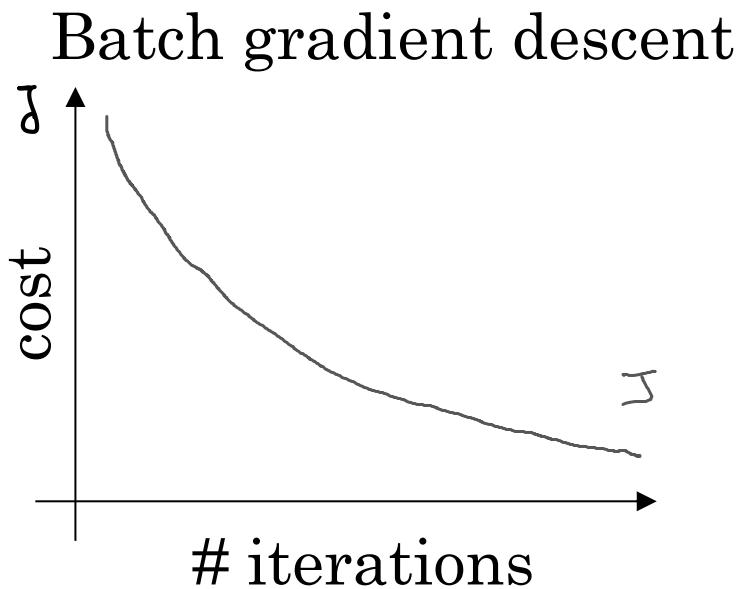
deeplearning.ai

Optimization Algorithms

Understanding mini-batch gradient descent

Training with mini batch gradient descent

کوئی خواہ سارے



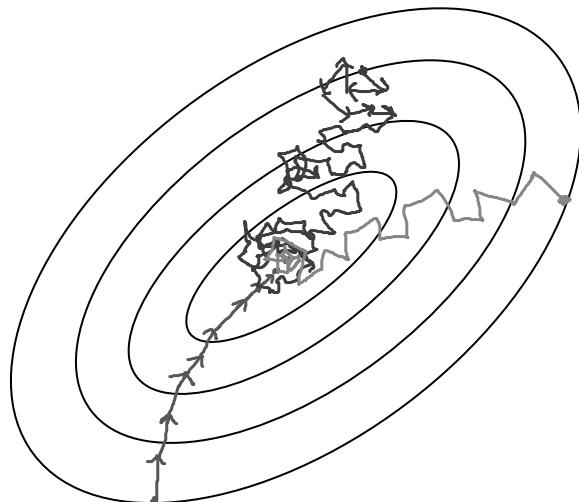
If mini-batch size = m = Batch gradient descent
 If mini-batch size = 1 = Stochastic gradient descent \Rightarrow mini-batch size between 1 and m

Choosing your mini-batch size

\rightarrow If mini-batch size = m : Batch gradient descent. $(X^{\{1\}}, Y^{\{1\}}) = (X, Y)$.

\rightarrow If mini-batch size = 1 : Stochastic gradient descent. Every example is its own
 $(X^{\{1\}}, Y^{\{1\}}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$ mini-batch.

In practice: Somewhere in-between 1 and m



Stochastic
gradient
descent

{
Use speedups
from vectorization

In-between
(mini-batch size
not too big/small)

Fastest learning.

- Vectorization.
 $(n \times 1000)$
- Make passes without
processing entire training set.

Batch
gradient descent
(mini-batch size = m)

↓
Two long
per iteration

Choosing your mini-batch size

If small train set : Use batch gradient descent.
($m \leq 2000$)

Typical mini-batch sizes:

$$\rightarrow 64, 128, 256, 512$$
$$2^6 \quad 2^7 \quad 2^8 \quad 2^9$$

$$\frac{1024}{2^{10}}$$

Make sure minibatch fits in CPU/GPU memory.

$$X^{[t]}, Y^{[t]}$$



deeplearning.ai

Optimization Algorithms

Exponentially weighted averages

هدف: یافتن روشی که سریعتر از gradient descent باشد

Temperature in London

$$\theta_1 = 40^{\circ}\text{F} \quad 4^{\circ}\text{C} \quad \leftarrow$$

$$\theta_2 = 49^{\circ}\text{F} \quad 9^{\circ}\text{C}$$

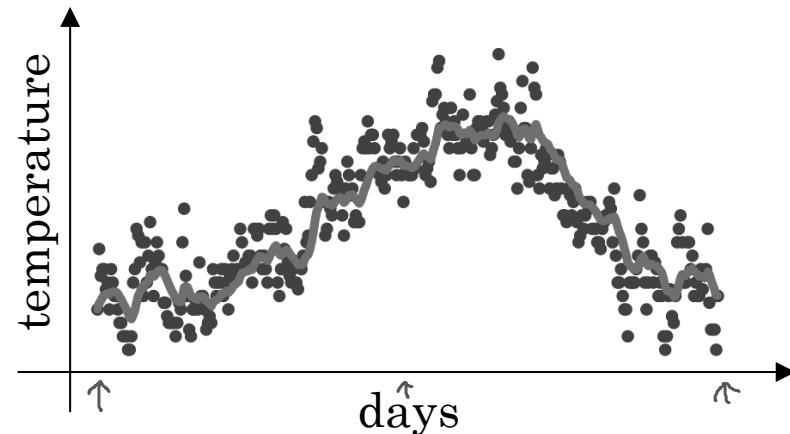
$$\theta_3 = 45^{\circ}\text{F} \quad \vdots$$

⋮

$$\theta_{180} = 60^{\circ}\text{F} \quad 15^{\circ}\text{C}$$

$$\theta_{181} = 56^{\circ}\text{F} \quad \vdots$$

⋮



$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

$$V_2 = 0.9 V_1 + 0.1 \theta_2$$

$$V_3 = 0.9 V_2 + 0.1 \theta_3$$

⋮

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

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t \Rightarrow \frac{1}{\beta-1} \text{ Previous } \theta$$

exponentially weighted moving average
exponentially weight averages

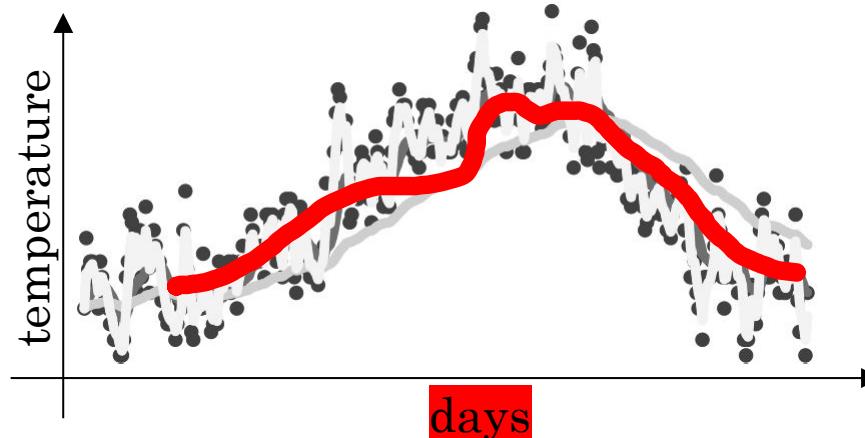
Exponentially weighted averages

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

$\underline{\beta} = 0.9$: ≈ 10 days' temperatur.
 $\underline{\beta} = 0.98$: ≈ 50 days
 $\underline{\beta} = 0.5$: ≈ 2 days

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

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





deeplearning.ai

Optimization Algorithms

Understanding exponentially weighted averages

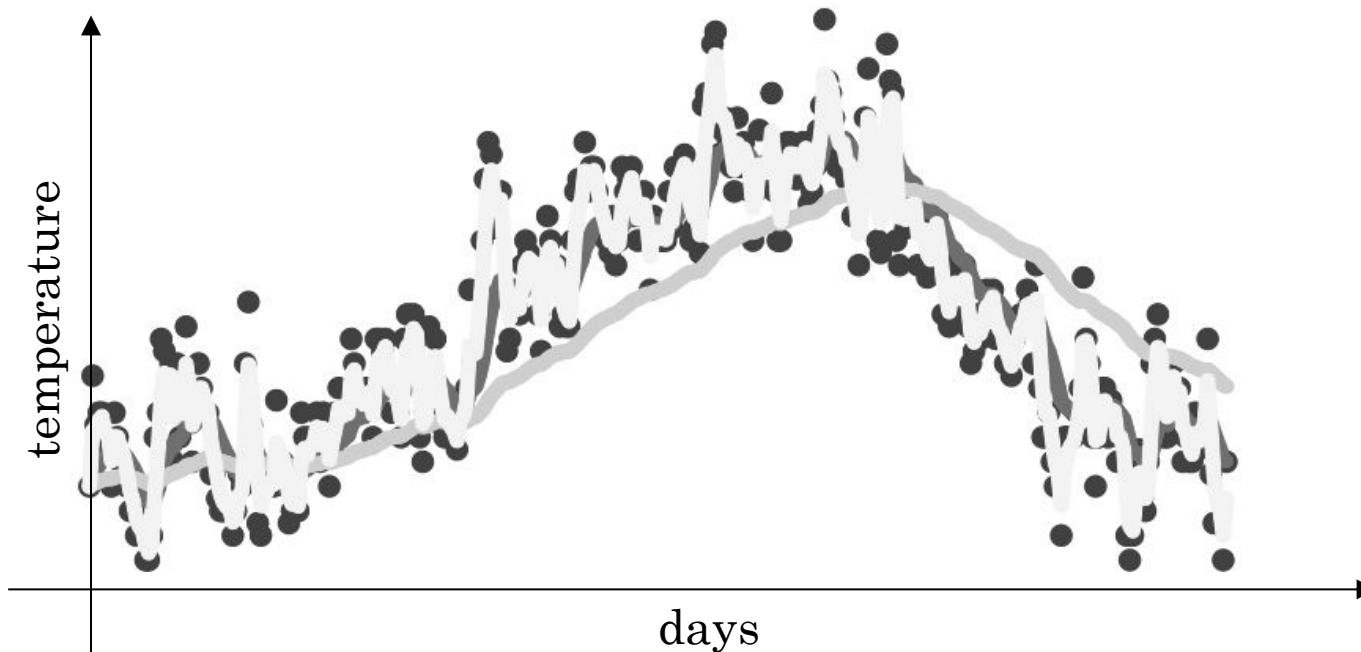
Exponentially weighted averages

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

$$\beta = 0.9$$

$$0.98$$

$$0.5$$



Exponentially weighted averages

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

$$v_{100} = 0.9v_{99} + 0.1\theta_{100}$$

$$v_{99} = 0.9v_{98} + 0.1\theta_{99}$$

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

...

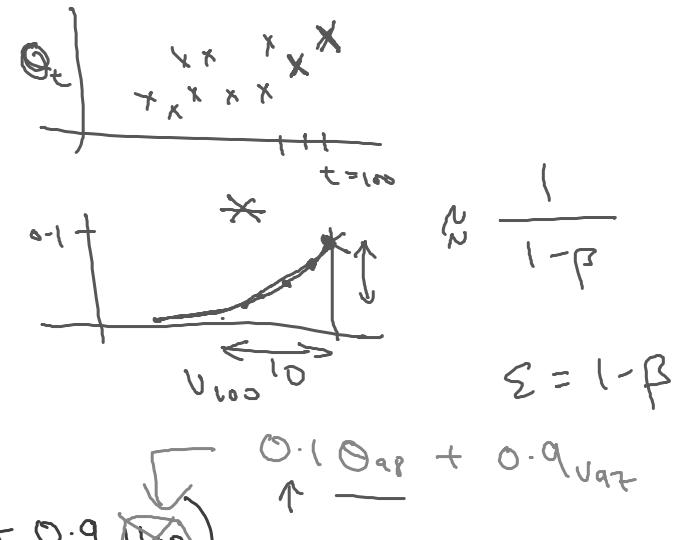
$$\underline{v_{100}} = 0.1 \theta_{100} + 0.9 \cancel{(0.1 \theta_{99})} + 0.9 \cancel{(0.1 \theta_{98})}$$

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

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

$$\frac{(1-\varepsilon)^{1/\varepsilon}}{0.9} = \frac{1}{e} \quad 0.98 ?$$

$\varepsilon = 0.02 \rightarrow \underline{0.98^{50}} \approx \frac{1}{e}$



Implementing exponentially weighted averages

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$

$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$v_3 = \beta v_2 + (1 - \beta) \theta_3$$

...

$$V_0 := 0$$

$$V_0 := \beta V + (1-\beta) \theta_1$$

$$V_0 := \beta V + (1-\beta) \theta_2$$

:

$$\rightarrow V_0 = 0$$

Repeat {

Get next θ_t

$$V_0 := \beta V_0 + (1-\beta) \theta_t \leftarrow$$

}

به صنایع بدای محاسب می‌اندیش است رون با خطای

Algorithm



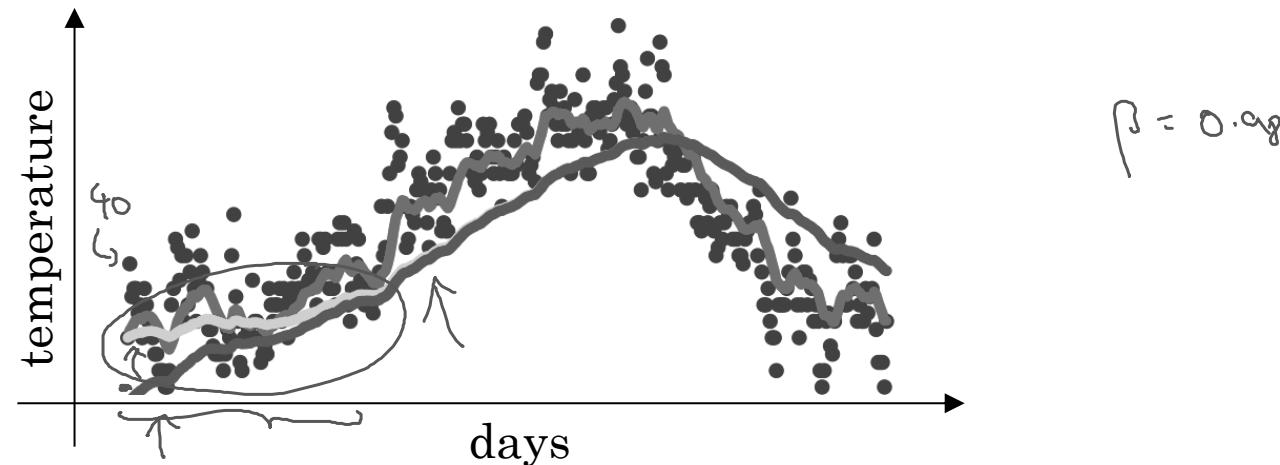
deeplearning.ai

Optimization Algorithms

Bias correction
in exponentially
weighted average

$$v_t = \frac{\beta v_{t-1} + (1-\beta) \theta_t}{1 - \beta^t} \Rightarrow \text{If } t \gg 1: 1 - \beta^t \approx 1 \Leftrightarrow \text{regularize کوچک کردن کوتاه مدت} \rightarrow \beta \text{ را کم کردن} \rightarrow \text{Bias Correction}$$

Bias correction



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

$$v_0 = 0$$

$$v_1 = 0.98 v_0 + 0.02 \theta_1$$

$$\begin{aligned} v_2 &= 0.98 v_1 + 0.02 \theta_2 \\ &= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2 \\ &= 0.0196 \theta_1 + 0.02 \theta_2 \end{aligned}$$

$$\left| \begin{array}{l} \frac{v_t}{1 - \beta^t} \\ t=2: 1 - \beta^t = 1 - (0.98)^2 = 0.0396 \\ \frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396} \end{array} \right.$$



deeplearning.ai

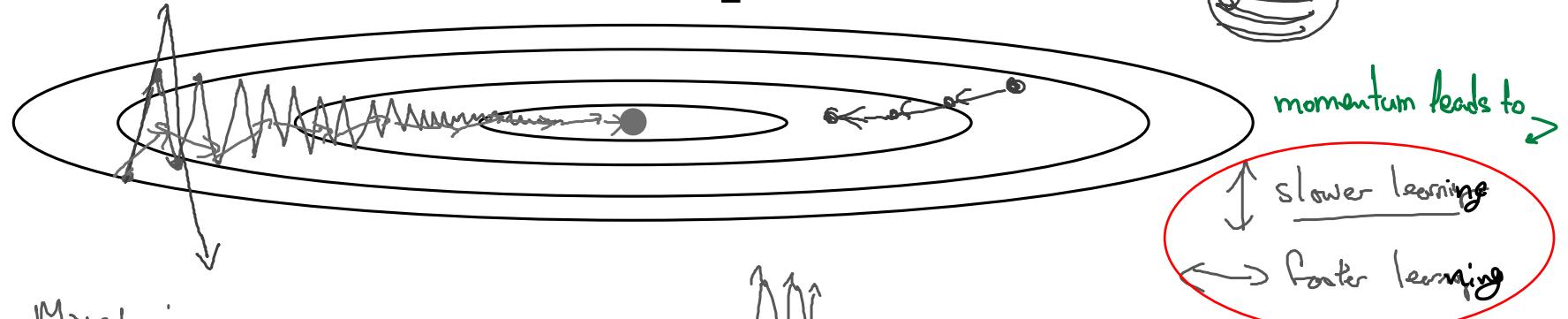
Optimization Algorithms

Gradient descent with momentum

Gradient descent with momentum

Ideas using exponentially weighted averages in gradient descent \Rightarrow always works better than gradient descent!

Gradient descent example



Momentum:

On iteration t :

Compute $d\omega, db$ on current mini-batch.

$$V_{d\omega} = \beta V_{d\omega} + (1-\beta) \cancel{d\omega}$$

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

Friction \uparrow velocity \uparrow acceleration

$$\theta_t = \theta_{t-1} + V_{\theta}$$

$$\omega = \omega - \alpha V_{d\omega}, \quad b := b - \alpha V_{db}$$

momentum leads to

\uparrow slower learning
 \downarrow faster learning

momentum (algorithm)

Implementation details

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

On iteration t :

Compute dW, db on the current mini-batch

$$\rightarrow v_{dw} = \beta v_{dw} + (1 - \beta) dW$$

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

$$W = W - \alpha v_{dw}, b = b - \alpha v_{db}$$

$$v_{dw} = \beta v_{dw} + dW \leftarrow$$

$$\cancel{v_{dw}} \rightarrow$$

$$\cancel{1 - \beta^t}$$

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

Hyperparameters: α, β
• α Hyperparam β

$$\beta = 0.9$$

average over last ≈ 10 gradients



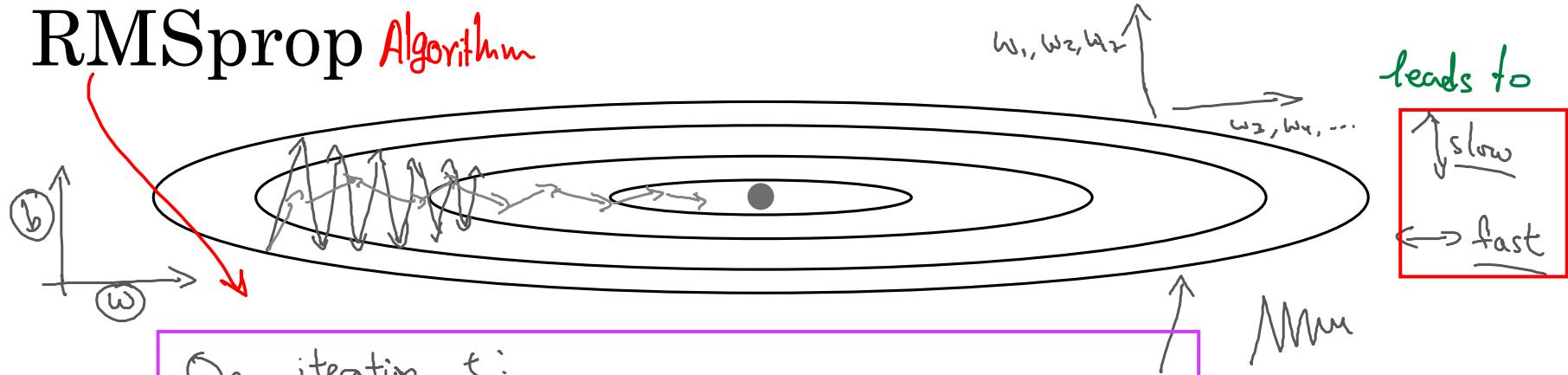
deeplearning.ai

Optimization Algorithms

RMSprop

Root mean Square Prop!

RMSprop Algorithm



On iteration t :

Compute dW, db on current mini-batch
element-wise

$$S_{dW} = \beta_2 S_{dW} + (1-\beta_2) dW^2 \quad \leftarrow \text{small}$$

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

$$w := w - \frac{\alpha}{\sqrt{S_{dW} + \epsilon}} dW \quad \leftarrow$$

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

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

\Rightarrow Possible to use $\alpha \gg 1$ without any diverging!



deeplearning.ai

Optimization Algorithms

Adam optimization algorithm

ئىچىرىنىڭدىن Adam, momentum, RMSprop ىلىكىن

$y\hat{=} \text{np.array}([.9, 0.2, 0.1, .4, .9])$

Adam optimization algorithm Algorithm

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

On iteration t :

Compute dW, db using current mini-batch

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

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

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

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

$$w := w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon}$$

$$b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

obj 4

Hyperparameters choice:

- learning rate α : needs to be tune
- $\rightarrow \beta_1 : 0.9 \rightarrow (\underline{dw})$
- $\rightarrow \beta_2 : 0.999 \rightarrow (\underline{dw^2})$
- $\rightarrow \epsilon : 10^{-8}$

Adam : Adaptive moment estimation

Adaptive moment estimation



deeplearning.ai

Optimization Algorithms

Learning rate decay

لیست α را نسبت به زمان t رام کم کنی.

Learning rate decay

Slowly reduce α



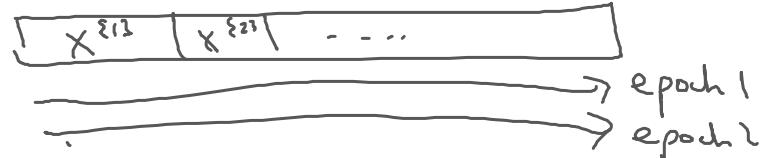
Learning rate decay

1 epoch = 1 pass through data.

$$\alpha = \frac{1}{1 + \underbrace{\text{decay-rate} * \text{epoch-num}}_{\downarrow} \alpha_0}$$

Epoch	α
1	0.1
2	0.067
3	0.05
4	0.04
:	:

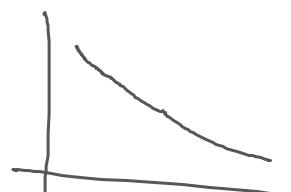
$\alpha \downarrow$



$$\alpha_0 = 0.2$$

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

learning rate decay



initial α

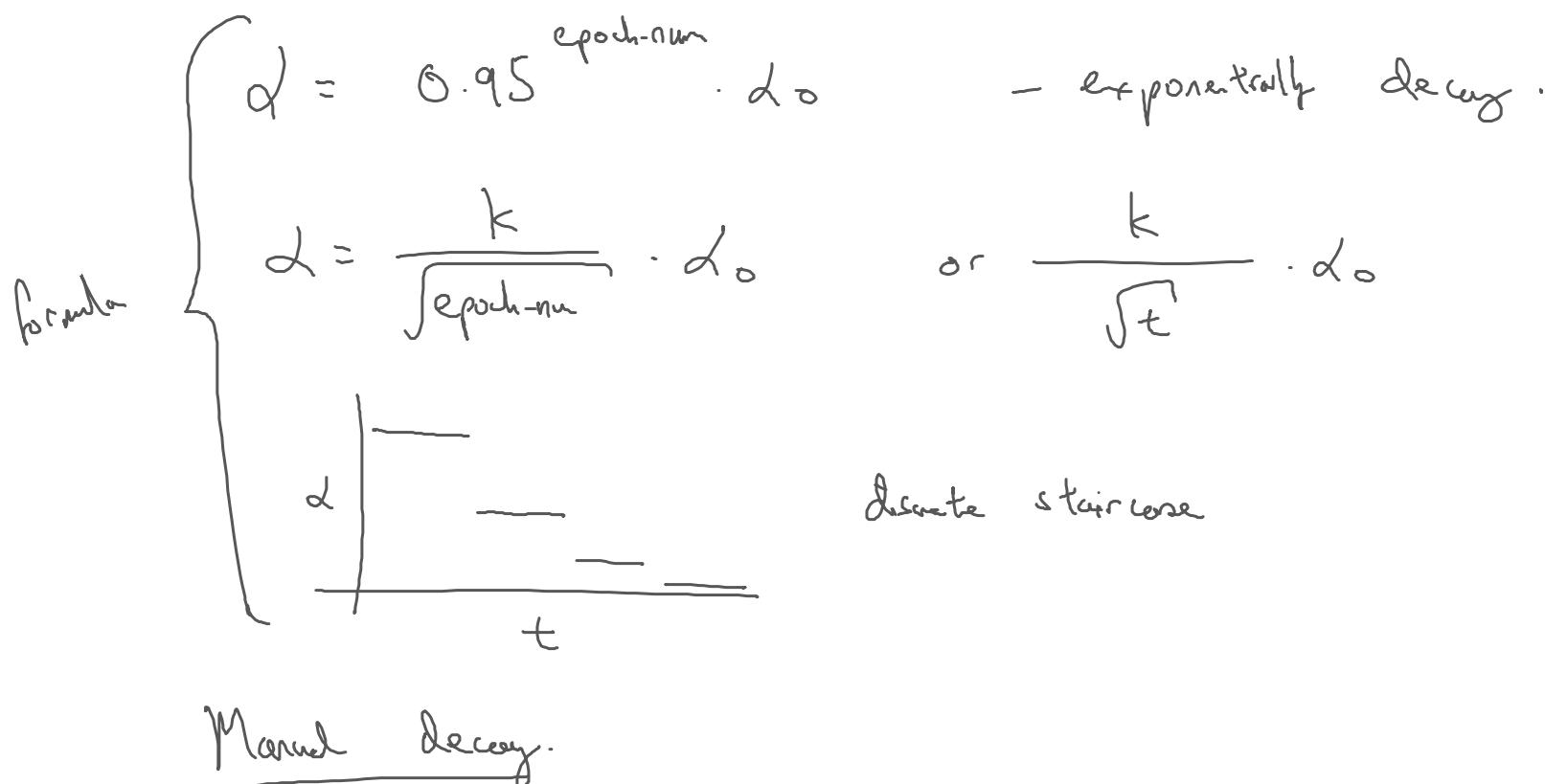
$$\alpha = \frac{\alpha_0}{1 + \underbrace{\text{decay-rate} * \text{epoch-num}}_{\text{Hyper Parameter}}}$$

سبیل
حکم

$$\alpha(t) = \frac{\alpha_0}{1 + \alpha t}$$

جدای از رویی قبل، کارهای دیگر هم ممکن است \Rightarrow چاچن است تابع $\alpha(t)$ نزولی تغیر گیند!

Other learning rate decay methods





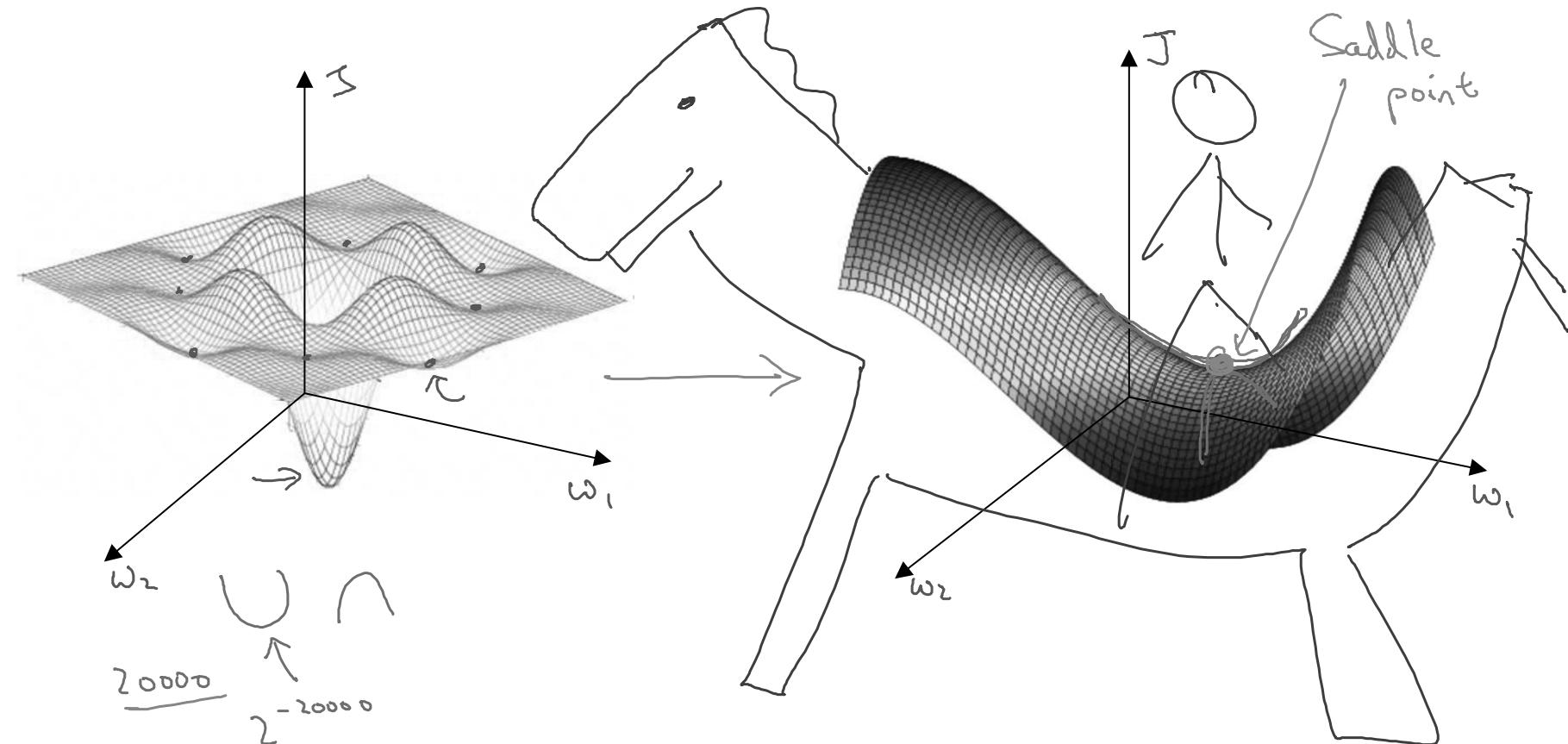
deeplearning.ai

Optimization Algorithms

The problem of local optima

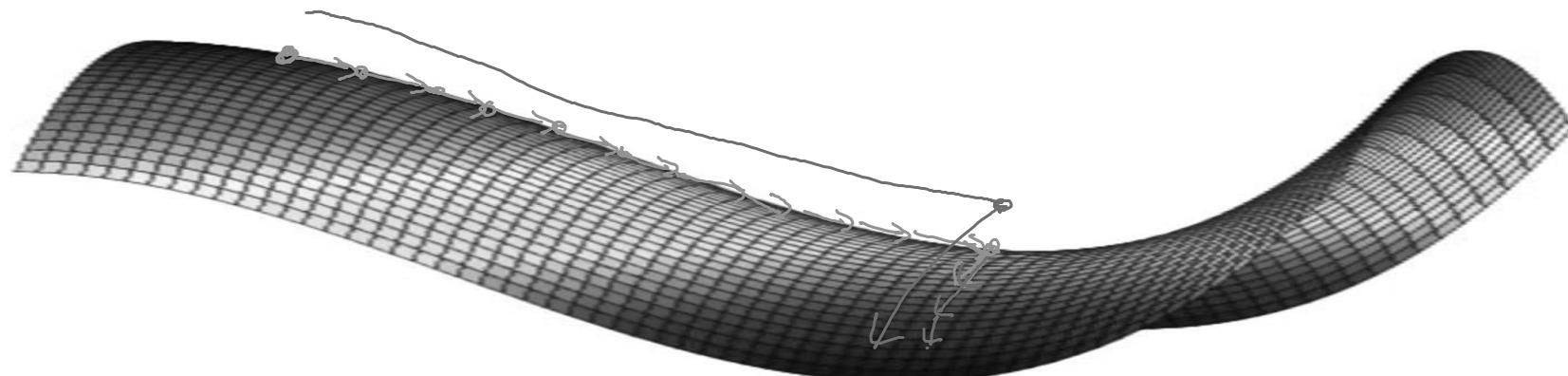
اگر نتایجی داشتیم که در آن صفر دارند، local min/max یا saddle points هستند. چون انتقام اینند در فضای n بعدی هم ممکن است جیز هم عالمت شوند. چنین پایانی است.

Local optima in neural networks



Problem of plateaus

فلاٹ



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow