

⇒ Mini-batch gradient descent

Vectorization allows you to efficiently compute on  $m$  examples.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$$

$(n \times m)$        $\underbrace{\quad}_{x^{(1)} \quad \dots \quad x^{(m)}} \quad x^{(1000)}$

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

min-batches of 1,000 each

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

$(1 \times m)$        $\underbrace{\quad}_{y^{(1)} \quad \dots \quad y^{(m)}} \quad y^{(5000)}$

$x^{(i)} \rightarrow$  training example

$z^{(L)} \rightarrow$  layer

$x^{(t)}, y^{(t)} \rightarrow$  mini-batch

Batch  $\Rightarrow$  entire training set some time

mini-batch  $\Rightarrow$  single mini batch,  $x^{(t)}, y^{(t)}$

Mini-batch gradient descent

for  $t = 1, \dots, 5000$

Forward prop on  $x^{(t)}$

$$z^{(1)} = w^{(0)} x^{(t)} + b^{(0)}$$

$$A^{(1)} = g^{(1)}(z^{(1)})$$

$\vdots$

$$A^{(L)} = g^{(L)}(z^{(L)})$$

Vectorized implementation

$\swarrow$  from  $x^{(t)}, y^{(t)}$

$$\text{Compute cost } J = \frac{1}{1000} \sum_{i=1}^L \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \underbrace{\frac{\lambda}{2 \cdot 1000} \sum_e \|w^{(e)}\|_F^2}_{\text{regularization}}$$

Backprop to compute gradients

$(dw, db)$

$J^{(t)}$

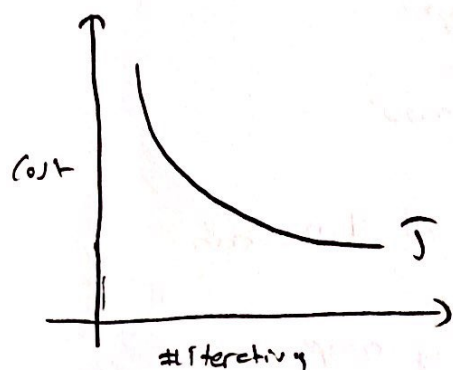
using  $(x^{(t)}, y^{(t)})$

$$w^{(e)} := w^{(e)} - \lambda dw^{(e)}, \quad b^{(e)} := b^{(e)} - \lambda db^{(e)}$$

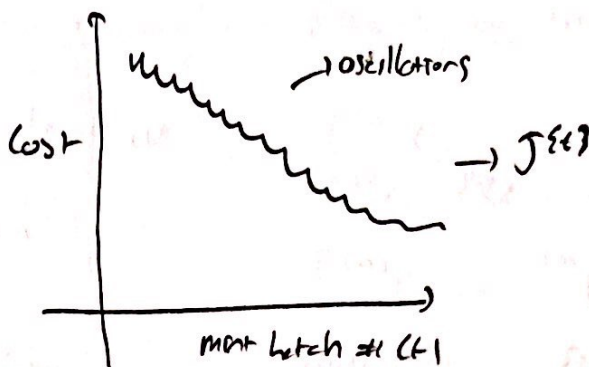
"1 epoch"  $\rightarrow$  pass through training set

# Understanding mini-batch gradient descent

Batch gradient descent



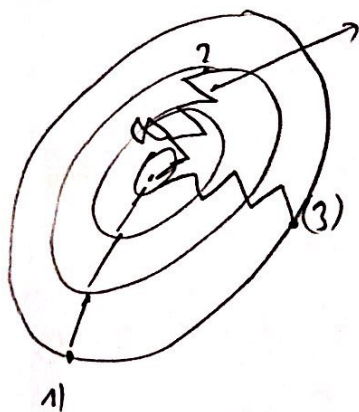
Mini-batch gradient descent



Choosing your mini-batch size #

1) If mini-batch size =  $m$  : Batch gradient descent.  $(x^{(1)}, y^{(1)}) = (x, y)$

2) If mini-batch size = 1 : Stochastic gradient descent. Every example is its own mini-batch.  
 $(x^{(1)}, y^{(1)}) = (x^{(1)}, y^{(1)})$



In practice somewhere in-between 1 and  $m$

Batch gradient descent

↓  
Too long per iteration

Stochastic gradient descent

↓  
lose speed up from vectors

Choosing your mini-batch size

• If small training set : use batch gradient des.  
 $(m \leq 2000)$

• If bigger training set : Typical mini-batch size

→ (64, 128, 256, 512) 2<sup>n</sup>  
 Common 1024 16384

Make sure minibatch fits in CPU/GPU memory  
 $x^{(k)}, y^{(k)}$

In-between

(mini-batch size)

Not big not small

↓  
Fastest learning

• vectorization

• Make progress without  
 entire training set

## Exponentially weighted averages

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

$$\beta = 0.9 \quad \approx 10 \text{ days ten}$$

$$\beta = 0.98 \quad \approx 50 \text{ days ten}$$

$$\beta = 0.9 \quad \approx 2 \text{ days ten}$$

$$\frac{1}{1-\beta} \text{ days temperature}$$

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

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

$$V_{100} = 0.1 \theta_{100} + 0.9 (V_{99} = 0.1 \theta_{99} + 0.9 (V_{98} = 0.1 \theta_{98} + 0.9 V_{97} \dots$$

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

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

$$\epsilon = 0.1 \text{ in this example}$$

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

$$\epsilon = 1-\beta$$

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

## Implementing exponentially weighted averages

$$V_0 := 0$$

$$V_t := \beta V_t + (1-\beta) \theta_t$$

$$V_t := \beta V_t + (1-\beta) \theta_t$$

$$V_0 := 0$$

Repeat

Get next  $\theta_t$

$$V_t := \beta V_t + (1-\beta) \theta_t$$



Beh:

correction in exponentially weighted average

$$\beta = 0.98$$

$$I = 0.98 V_0 + 0.02 Q_1$$

$$I_2 = 0.98 I_1 + 0.02 Q_2$$

Vo' o olday ica ille kismiler  
notch checked!!

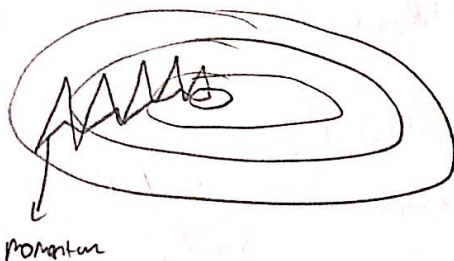
Bias correction ile  $V_0 = 0$  same corol

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

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

$$\frac{V_2}{0.0396} = \frac{0.0146 Q_1 + 0.02 Q_2}{0.0396}$$

Gradient descent with momentum



Momentum

on iteration  $t$ :

compute  $\Delta w$  &  $\Delta b$  on current mini-batch

$$V \Delta w = \beta V \Delta w + (1 - \beta) \Delta w$$

(other omitted)

$$V \Delta b = \beta V \Delta b + (1 - \beta) \Delta b$$

(how kullu)

$$w = w - \eta V \Delta w, b = b - \eta V \Delta b$$

$$V_0 = \beta V_0 + (1 - \beta) Q_1$$

$$\frac{V \Delta w}{1 - \beta^t} \Rightarrow \text{bias correction kullu}$$

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

RMSprop (Root Mean Squared)

~~ADA~~

↓ sbu } Bunun için RMSprop kullanılır.  
↔ fet }

On iteration t:

Compute  $d_w, d_b$  on current mini-batch

$$S_{dw} = \beta S_{dw} + (1-\beta) d_w^2 \rightarrow \text{element wise}$$

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

$$w := w - \alpha \frac{d_w}{\sqrt{S_{dw} + \epsilon}} \quad b := b - \alpha \frac{d_b}{\sqrt{S_{db} + \epsilon}}$$

↗  $\epsilon \rightarrow 10^{-7}$   $\checkmark$   $\rightarrow$   $\epsilon$  değeri çok küçük  
sıfıra yakın değere

↗  $y_i$   $\rightarrow$   $y_i$  değeri  
↘  $w$

Adam optimization algorithm

Momentum + RMSprop

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

On iteration t:

Compute  $d_w, d_b$  using current mini batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) d_w, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) d_b \rightarrow \text{momentum}$$

$$S_{dw} = \beta_2 V_{dw} + (1-\beta_2) d_w^2 \quad S_{db} = \beta_2 V_{db} + (1-\beta_2) d_b^2 \rightarrow$$

biçim  $\rightarrow$   $\epsilon$  değeri çok küçük

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

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

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

$\epsilon \rightarrow 10^{-7}$

Hyperparameters choice

Adam = Adaptive Moment Estimation

$\alpha$  needs to be fine

$\beta_1 = 0.9$  ( $d\omega$ )  $\rightarrow$  First moment

$\beta_2 = 0.999$  ( $d\omega^2$ )  $\rightarrow$  Second moment

$\epsilon = 10^{-7}$  (not important)

Learning rate decay

Hedete globalstika kiziltnesiye gide

1 epoch = (pass) through data

$$\alpha = \frac{1}{1 + \text{decay-rate} \times \text{epoch-num}} \cdot \alpha_0$$

Epoch	$\alpha$	$\alpha_0 = 0.2$ decay rate = 1
1	0.1	
2	0.67	
3	0.5	
4	0.4	

Hedete globalstika kiziltnesiye gide  
kiziltnesiye minimumu

Other learning rate decay methods #

1)  $\alpha = 0.55^{\text{epoch-num}} \cdot \alpha_0 \rightarrow$  exponentially decay

2)  $\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0$  or  $\frac{k}{\sqrt{t}} \cdot \alpha_0$

Menel decay  
ele indirgen.

