

§ Improving DNN — week 2 Optimization Algorithms

1. Batch VS. mini-batch

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} & x^{(1001)} & \dots & x^{(2000)} & \dots & x^{(m)} \end{bmatrix}$$

(n_x, m) $X^{(1)}$ $X^{(2)}$ $X^{(5000)}$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} & y^{(1001)} & \dots & y^{(2000)} & \dots & y^{(m)} \end{bmatrix}$$

$(1, m)$ $Y^{(1)}$ $Y^{(2)}$ $Y^{(5000)}$

for $t = 1, \dots, 5000$

forward prop on $X^{(t)}$

$$Z^{[1]} = W^{[1]} X^{(t)} + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$\vdots$$

$$A^{[L]} = g^{[L]}(Z^{[L]})$$

1000 examples

1 stop of gradient descent using $X^{(t)}, Y^{(t)}$ (as if $m = 1000$)

from $X^{(t)}, Y^{(t)}$

compute cost $J = \frac{1}{1000} \sum_{i=1}^1 \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_f \|W^{[L]}\|_F^2$

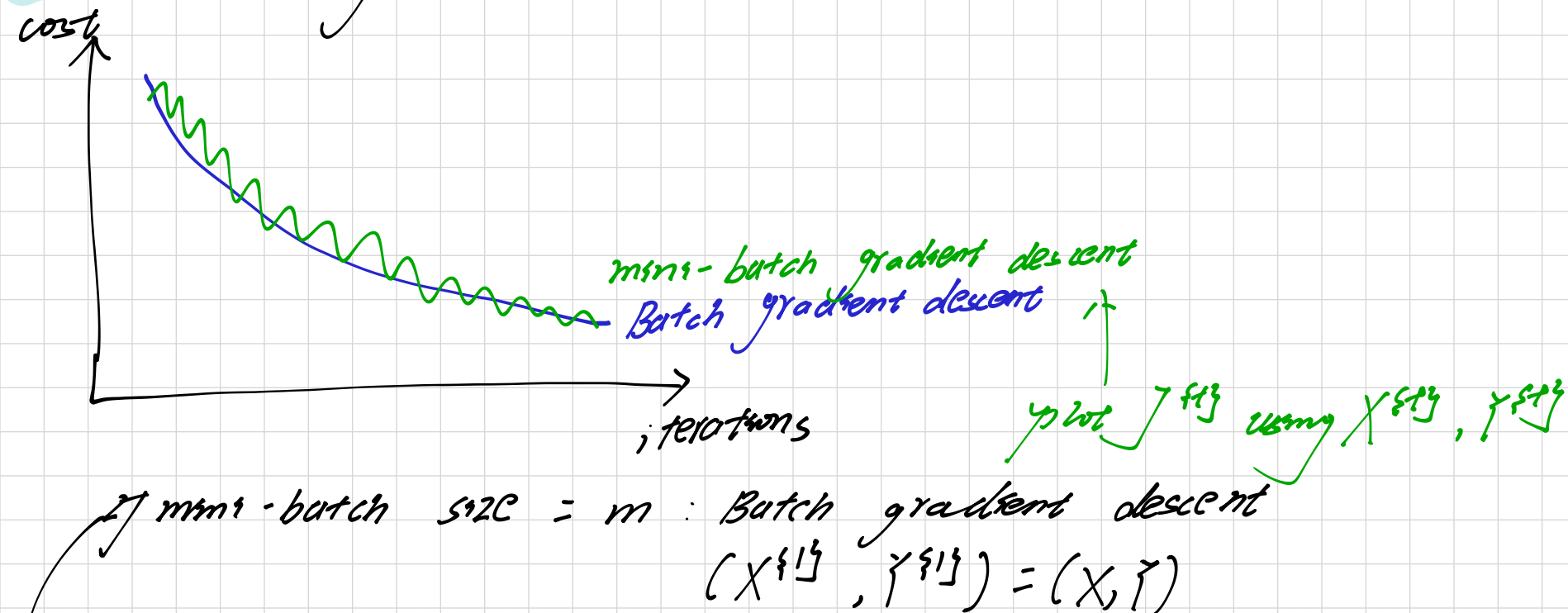
Backprop to compute cost $J^{(t)}$ (using $X^{(t)}, Y^{(t)}$)

$$W^{[L]} := W^{[L]} - \alpha dW^{[L]}$$

$$b^{[L]} := b^{[L]} - \alpha db^{[L]}$$

1 epoch → passing through training set

2. mini-batch gradient descent



2. mini-batch size = 1 : Stochastic gradient descent

$$(x^{(1)}, y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$$

In practice : between 1 and m

Batch gradient descent

too long per iteration

In between

faster learning

Stochastic gradient descent

low vectorization

{ small training set : use batch gradient descent
($m \leq 2000$)

Typical mini-batch size : 64, 128, 256, 512, 1024

* Make sure mini-batch fit in CPU / GPU memory

3. Exponentially weighted averages

$$v_0 = 0$$

$$v_\theta := 0$$

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

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

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

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

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

$$\vdots$$

...

$$v_\theta := 0$$

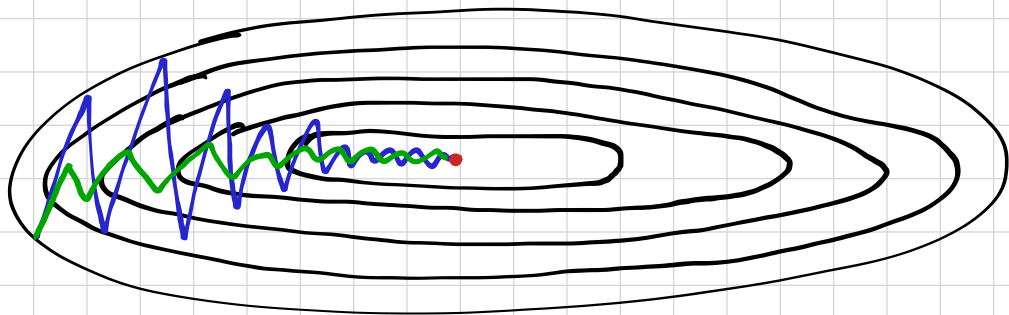
Repeat {

Get next θ_t

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

}

4. Gradient Descent with momentum



if wants \downarrow slower learning

\longleftrightarrow faster learning

\Rightarrow use momentum

Momentum:

On iteration t :

Compute dw, db on correct mini-batch

$$\begin{aligned} V_{dw} &= \beta dw + (1-\beta) V_{dw} \\ V_{db} &= \beta db + (1-\beta) V_{db} \end{aligned} \quad \left\} \Leftarrow V_0 = \beta V_0 + (1-\beta) \theta_t$$

friction $W := W - \alpha dw$ acceleration
 $b := b - \alpha db$ velocity

5. RMSprop

On iteration t :

Compute dw, db on mini-batch

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

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

$$W := W - \alpha \frac{dw}{\sqrt{S_{dw}} + \epsilon} \quad \leftarrow \epsilon = 10^{-8}$$

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

6. Adam optimization algorithm

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

On iteration t :

Compute dw, db using current mini-batch

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

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

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

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

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

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

$$W := W - \alpha \frac{V_{dw}^{current}}{\sqrt{S_{dw}} + \epsilon}$$

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

hyperparameters choice:

α : need to be tune

β_1 : $0.9 \rightarrow (dw)$

β_2 : $0.999 \rightarrow (dw^2)$

ϵ : 10^{-8}