# Mini - Batch Gradient Descent

e.j.  $m = 1,000,000$

$\rightarrow$ split into 1000 mini-batches of 1,000 each

Mini-Batch $t$ :  $X^{\{t\}}$, $Y^{\{t\}}$
$\quad\quad\quad\quad\quad\quad\quad\quad (n_x, 1000) \quad (1, 1000)$

for $t = 1$ to 1000 :

Forward Prop on $X^{\{t\}}$

$$Z^{[1]} = W^{[1]} X^{\{t\}} + b^{[1]}$$
$$A^{[1]} = g^{[1]}(Z^{[1]})$$
$$\vdots$$
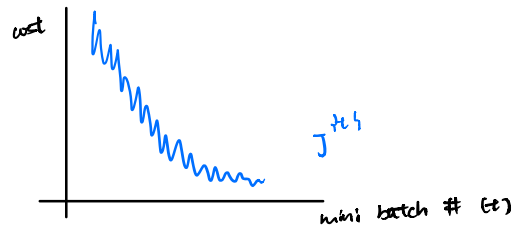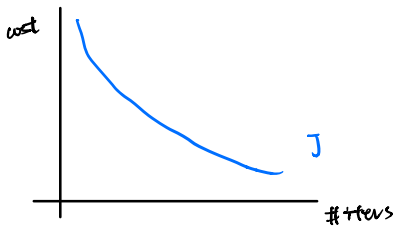$$A^{[L]} = g^{[L]}(Z^{[L]})$$

$\quad\quad\quad\quad$ } Vectorized implementation

Compute $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{\ell} \| w^{[\ell]} \|_F^2$

Backprop to compute gradient w.r.t. $J^{\{t\}}$ ( using $(X^{\{t\}}, Y^{\{t\}})$ )

$$W^{[\ell]} := W^{[\ell]} - \alpha \, dW^{[\ell]} \quad ; \quad b^{[\ell]} := b^{[\ell]} - \alpha \, b^{[\ell]}$$

'1 epoch' — pass through training set



$\quad$ cost $\quad\quad$ J $\quad\quad$ # iters

$\quad$ cost $\quad\quad$ $J^{\{t\}}$ $\quad\quad$ mini batch # (t)

Batch Gradient Descent $\xleftarrow{\text{size} = m}$ Mini Gradient Descent

$\quad\quad$ stochastic $\sim$ $\swarrow$ size = 1

$\quad\quad\quad$ ( lose speedup from vectorisation )

$\Rightarrow$ typical mini-batch size :

$\quad\quad\quad$ 64, 128, 256, $\cdots$

make sure mini-batch fit in CPU/GPU memory

## Exponentially Weighted Averages

e.g.
$$\theta_1 = 40°F$$
$$\theta_2 = 49°F$$
$$\vdots$$

$$\Rightarrow \quad V_0 = 0$$
$$V_1 = 0.9 V_0 + 0.1 \theta_1$$
$$V_2 = 0.9 V_1 + 0.1 \theta_2$$
$$\vdots$$
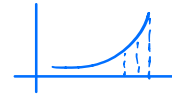$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

e.g.
$$V_{100} = 0.1 \theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1 \times (0.9)^2 \theta_{98} + \cdots$$

$$0.9^{10} \approx 0.35 \approx \frac{1}{e}$$
$$(1 - \varepsilon)^{1/\varepsilon} = \frac{1}{e}$$

after $\frac{1}{1-\beta}$ days:
about 1/3 left

$$\Rightarrow \quad V_\theta = 0$$
Repeat {

       Get next $\theta_t$

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

}

temperature



days

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

averages over $\approx \frac{1}{1-\beta}$ days
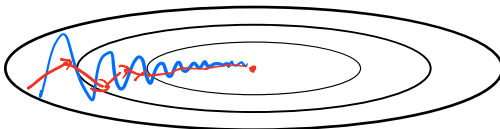


Problem:
$$V_1 = \beta V_0 + (1-\beta) \theta_1$$
$$= (1-\beta) \theta_1$$

Bias Correction $\Rightarrow \dfrac{V_t}{1 - \beta^t}$

(only matters in early stage)

## Gradient Descent w/ Momentum



$\updownarrow$ slower
$\leftrightarrow$ faster

Momentum:
On Iteration t:

    Compute $dw$, $db$ on current mini-batch

$$V_{dw} = \beta V_{dw} + (1-\beta) dw$$
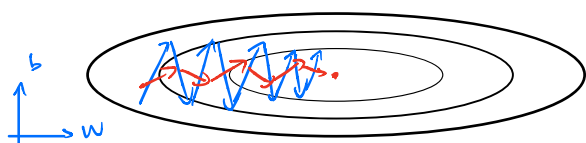$$V_{db} = \beta V_{db} + (1-\beta) db$$
$$w := w - \alpha V_{dw}, \quad b := b - \alpha V_{db}$$

another form:
$\beta V_{dw} + dw \quad \leftarrow$ tune $\alpha$

($\alpha, \beta$ hyperparameters)

## RMSProp



$\downarrow$ slower
$\leftrightarrow$ faster

On iteration $t$,

    compute $dw, db$ on current min-batch

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

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

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

Prevent zero error
e.g. $\epsilon = 10^{-8}$

## Adam Optimization Algorithm — Adaptive Moment Estimation

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

On iteration $t$,

    Compute $dw, db$ using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw \;,\; V_{db} = \beta_1 V_{db} + (1-\beta_1) db \quad \swarrow \text{Momentum}$$

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

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

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

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

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

hyperparameters:
    $\alpha$: needs to be tuned
    $\beta_1 : 0.9$     $(dw)$
    $\beta_2 : 0.999$     $(dw^2)$
    $\epsilon : 10^{-8}$

## Learning Rate Decay



Decay

1 epoch = 1 pass of data

$$\alpha = \frac{1}{1 + \text{decay rate} * \text{epoch \#}} \alpha_0$$

$\cdots 0.2$

$\Rightarrow$ e.g.

| epoch | 1 | 2 | 3 | 4 | ... |
|-------|-----|------|------|------|-----|
| $\alpha$ | 0.1 | 0.06 | 0.05 | 0.04 | ... |

Other decay methods:

$$\alpha = 0.95^{\text{epoch \#}} \, \alpha_0$$

$$\alpha = \frac{k}{\sqrt{\text{epoch \#}}} \, \alpha_0$$

$$\alpha = \quad$$



## Problem of Local Optima

In high dimension, saddle points are more common

Unlikely to stuck in a bad local optima

Plateaus can make learning slow