# Optimization Algorithms - Coursera

SGD
Momentum          } — Optimization methods.
RmsProp
Adam

→ Use Random mini batches to accelerate the convergence and improve the optimization.

## Batch Gradient Descent vs. Mini Batch Gradient Descent

Suppose records are $5000000$

$$\downarrow 5000 * 1000 \text{ each batches}$$

for each batch $t = 1, \ldots 5000$

{ Forward prop
  → Compute Activation fns → (Vectorized implementation on 1000 samples)

  $fn(x^{[t]}, y^{[t]})$

Compute Cost for

$$J^{(t)} = \frac{1}{1000} \sum_{i=1}^{L} L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \times 1000} \sum_{l} \|w^{[l]}\|^2$$

Back propagation to compute gradients wrt $J^{(t)}$ using $x^{[t]}, y^{[t]}$

$$w^{(l)} := w^{[l]} - \alpha \, dw^{[l]}, \quad b^{[l]} := b^{[l]} - \alpha \, db^{[l]}$$

}

1 pass thru your training set — Doing 1 'epoch' of training.
(using mini batch gradient descent)

'epoch' → One single pass thru the gradient descent. training set
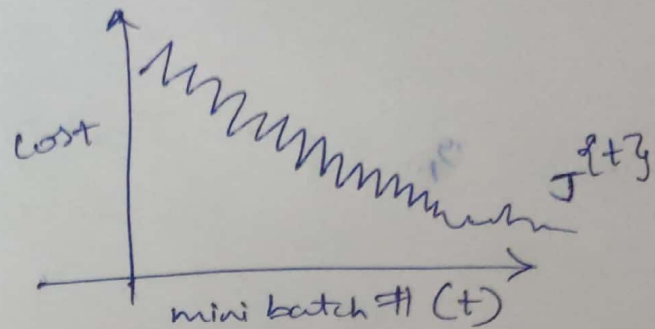
In Batch GD — 1 epoch allows us to take 1 gradient descent step

In mini Batch GD — 1 epoch allows us to take 5000 gradient descent steps.

You might want to have another for loop to include
~~wong~~ passes. Thru the training set. [i.e many epochs]
multiple

Batch Gradient Descent.          mini-batch gradient descent.



Cost
#iterations → J

cost
mini batch #1 (t)      $J^{\{t\}}$

Training process of both.

Choosing your ~~&~~ mini batch size ( See snapshot )

Typical mini-batch sizes:
            64, 128, 256 . . . . . $2^n$

Make sure all $x^{\{t\}}, y^{\{t\}}$ fits in CPU/GPU memory.

### Exponentially weighted averages.

        & also called

Exponentially weighted moving averages

                Moving average is calculated as,

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



temperature

days

$$\boxed{V_t = \beta V_{t-1} + (1-\beta)\theta_t}$$ → Formula to implement E.W. moving Avg

Here $\beta = 0.9$, $V_t$ is approximately
                average over $\approx \dfrac{1}{1-\beta}$ days temp.

So if.

$\beta = 0.9$ : $\approx$ 10 days
$\beta = 0.98$ : $\approx$ 50 days → more smooth the line (since we are averaging
                                            on much larger window)
$\beta = 0.5$ : $\approx$ 2 days → more noise

# Understand Exponentially weighted Averages.

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

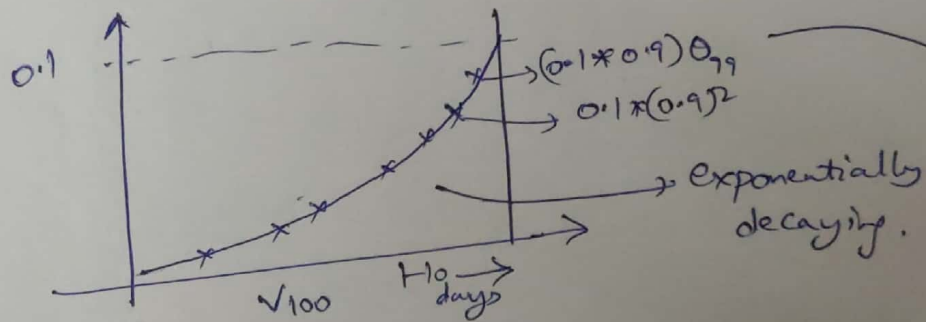$$V_{100} = 0.9 V_{99} + 0.1 \theta_{100}$$

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

$$\vdots$$

$$V_{100} = 0.1 \theta_{100} + 0.1 * 0.9 \, \theta_{99} + 0.1 * (0.9)^2 \theta_{98} + 0.1 * (0.9)^3 \theta_{97} + \cdots$$

One way to show this in pictures



$$0.1, \quad 0.1 * 0.9, \quad 0.1 * (0.9)^2 \cdots \longrightarrow \text{all These add up to } 1$$
upto a detail called
$$\boxed{\text{Bias correction}}.$$

??? (refer to previous page)

How many days is this averaging over? i.e $\left(\frac{1}{1-\beta}\right)$

→ its averaging over 10 days

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

So it takes 10 days for this height to decay over $\frac{1}{3}$rd.

****

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

# Bias Correction in exponentially weighted averages.



temperature

day

→ expertly.

→ actual
(starts off
really low)

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

$V_0 = 0$

$V_1 = 0.98 V_0 + 0.02 \theta_1$

$V_1 = 0.02 \theta_1$

$V_2 = 0.98 (0.02 \theta_1) + 0.02 \theta_2$ } → very slow start to the weighted averages.

to counter this, use $\dfrac{V_t}{1-\beta^t}$ as $V_t$

for example $t = 2$:

$$\frac{V_2}{1-(0.98)^2} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$ } → This will give it a better start.

# Gradient descent with Momentum.

The basic idea is to compute an exponentially weighted average of your gradients and use that gradient to update your weights instead.



↓ Slower learning

↔ faster Learning.

 → ball rolling down a bowl.

## Implementation details.

On iteration t :

(prevents from speeding up) Compute $dw, db$ on the current mini-batch.

friction $\quad V_{dw} = \beta V_{dw} + (1-\beta) dw \quad$ → acceleration

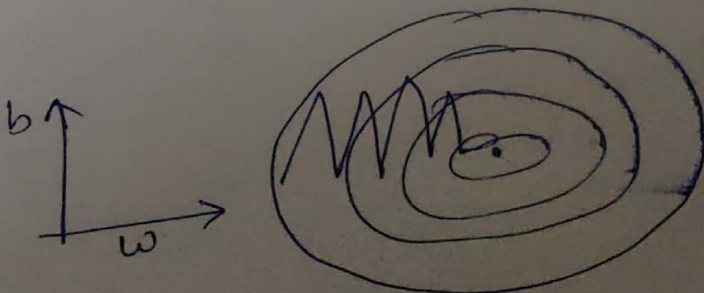$\qquad V_{db} = \beta V_{db} + (1-\beta) db$

velocity

$\qquad w = w - \alpha V_{dw}, \quad b = b - \alpha V_{db}$

Hyperparameters: $\alpha, \beta, \quad \underbrace{\beta = 0.9}_{\text{Common value we use}} \qquad$ iterations

i.e Average over the last $10 \land$ gradients (or iteration)

## RMSProp

Root Mean Square prop.



b ↑ \
↳ w →

↑ slow \
↓ b

↔ fast \
w

On iteration t:

Compute $dw, db$ on current mini-batch

$S_{dw} \pm \beta S_{dw} + (1-\beta) dw^2$    element-wise
                                   Squaring operation
                                         $\rightarrow$ Small

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

$$w := W - \alpha \frac{dw}{\sqrt{S_{dw}}} \quad , \quad b = b - \alpha \frac{db}{\sqrt{S_{db}}}$$

Updates too in horizontal direction is large.

So Updates in vertical direction is small

$\rightarrow$ This is keeping an exponentially weighted average of the square of the derivatives.

## Adam optimization Algorithm.

$\rightarrow$ Momentum + RMSProp

While implementing Adam, we will need to add bias correction.

for $V_{dw}, V_{db}, S_{dw}, S_{db}$.

finally,

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

Hyper parameters

$\alpha$: needs to be tuned

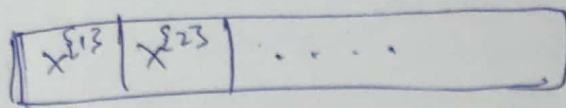$\beta_1 : 0.9 \rightarrow dw$
$\beta_2 : 0.99 \rightarrow dw^2$   Use these default values
$\epsilon : 10^{-8}$

[ Adam - Adaptive moment Estimation ]

# Learning Rate decay

↳ To slowly reduce learning rate over time.

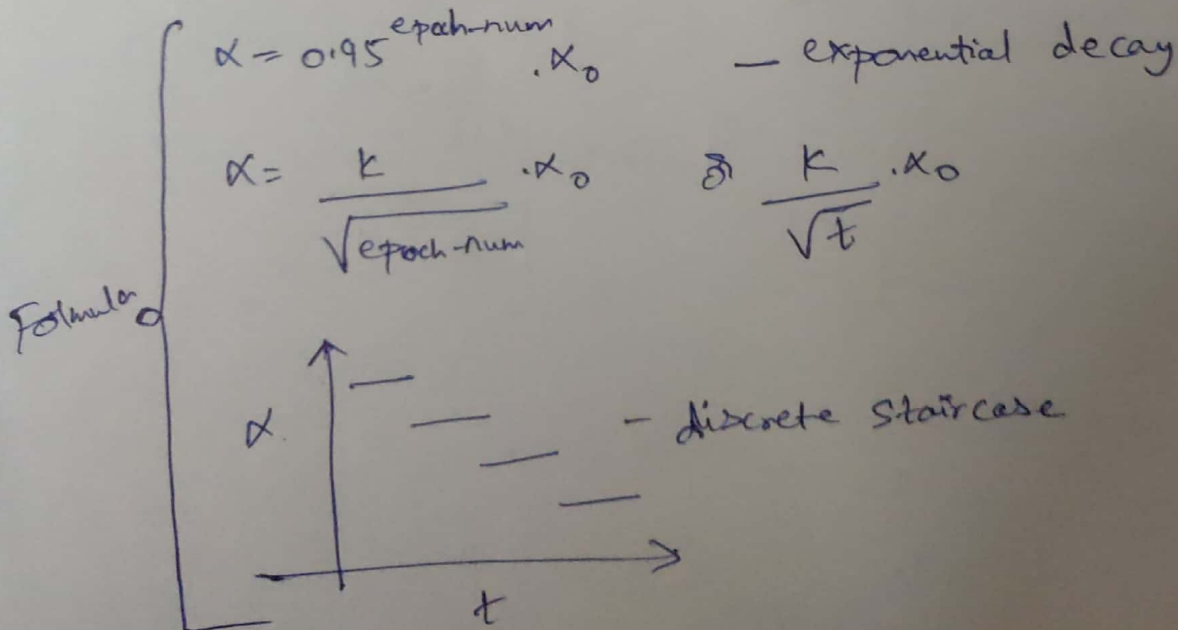$$[\; x^{\{1\}} \mid x^{\{2\}} \mid \ldots \ldots \;]$$

→ epoch 1

→ epoch 2

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

Decay rate is another hyperparameter to tune.

| Epoch | $\alpha$ |
|-------|----------|
| 1 | 0.01 |
| 2 | 0.067 |
| 3 | 0.05 |
| 4 | 0.04 (Decays) |

$\alpha_0 = 0.2$

decay rate = 1

Formula
$$\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \qquad - \text{exponential decay}$$

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



— discrete staircase

manual decay { watch your model as its training and if learning rate has slowed down, manually increase $\alpha$.