

Optimization algorithms

Mini-batch gradient descent (GD we had before is called "Batch Gradient Descent" > all data computed at once)

→ Vectorization allows us to efficiently compute on m examples

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \quad \begin{matrix} (n_{x,m}) \\ \rightarrow \text{can still be slow if } m \text{ is large} \end{matrix}$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix} \quad \begin{matrix} (1,m) \\ \rightarrow \text{What if } m = 5,000,000? \end{matrix}$$

→ mini-batch GD will perform much better on large data sets

→ Possible to get faster algorithm → split training set into smaller mini-batches each mini-batch has $\sim 1,000$ training examples

$$X = \left[\underbrace{x^{(1)} \dots x^{(1000)}}_{X^{(1)}} \mid \underbrace{x^{(1001)} \dots x^{(2000)}}_{X^{(2)}} \mid \dots \mid \underbrace{x^{(5000001)} \dots x^{(5000000)}}_{X^{(5000)}} \right]$$

new notation for mini-batches

→ $X^{(1)} \dots X^{(5000)}$ for $m = 5,000,000$ we get 5000 mini-batches of size 1000

same for $Y = \left[\underbrace{y^{(1)} \dots y^{(1000)}}_{Y^{(1)}} \mid \underbrace{y^{(1001)} \dots y^{(2000)}}_{Y^{(2)}} \mid \dots \mid \underbrace{y^{(5000001)} \dots y^{(5000000)}}_{Y^{(5000)}} \right]$

mini-batch t : $X^{(t)}, Y^{(t)}$

How to run?

1 step of GD using $X^{(t)}, Y^{(t)}$ (as if $m=1000$)

- for $t = 1 \dots 5000$
- Forward Prop on $X^{(t)}$

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

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

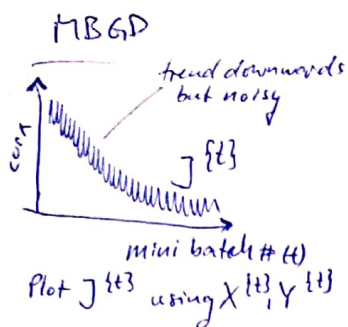
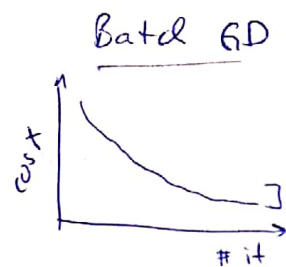
$$\vdots$$

$$A^{[L]} = g^{[L]}(z^{[L]})$$
 - Compute cost of one mini-batch
$$J^{(t)} = \frac{1}{1000} \sum_{i=1}^L \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2 \cdot 1000} \|W^{[2]}\|_F^2$$
 - Backpropagation to compute gradients w.r.t. $J^{(t)}(X^{(t)}, Y^{(t)})$
 - ↳ update weights: $W^{[2]} := W^{[2]} - \alpha dW^{[2]}, b^{[2]} := b^{[2]} - \alpha db^{[2]}$
- from $X^{(t)}, Y^{(t)}$

→ "1 epoch" of training → 1 pass through training set

Batch-GD: single pass through training set allows to only take 1 GD step
 Mini-batch-GD: ————— allows to take 5000 GD steps (in this example)

Understanding mini-batch gradient descent



oscillations from plot because $X^{(1)}, Y^{(1)}$ has easy cost value but $X^{(2)}, Y^{(2)}$ might have higher cost than 1^{st}

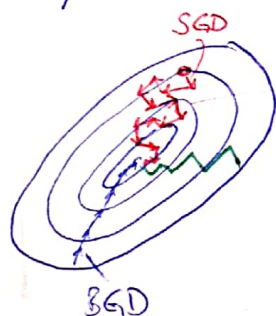
Parameter to choose

→ Choosing your mini-batch size!

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

extreme 2 - If mini-batch size = 1: Stochastic Gradient Descent $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$ Every example is its own mini-batch

In practice mini-batch size will be in between 1 and m → looking at one training example at a time



will never really converge

SGD

lose speed up from vectorization

In practice In-between

(mini-batch size not too high/small)

↓

Fastest learning

- can make use of vectorization
- make progress without needing to wait to process the entire training set

BGD

(mini-batch size = m)

↓

too long per iteration

Guidelines:

□ If small dataset: Use BGD ($m < 2000$)

□ typical mini-batch size:

64, 128, 256, 512

power of 2: $2^6, 2^7, 2^8, 2^9$

more common

□ Make sure mini-batch fits in CPU/GPU memory

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

→ mini-batch size is also a hyperparameter

Exponentially weighted averages

ex. temperature in London

$\theta_1 = 4^\circ\text{C}$

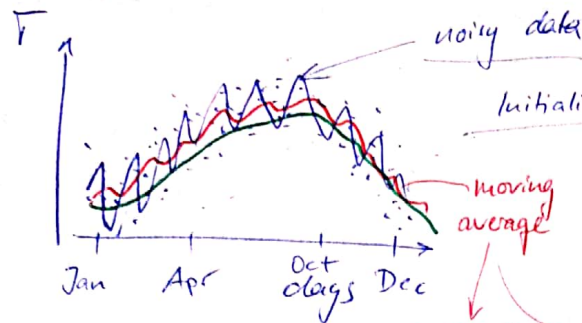
$\theta_2 = 9^\circ\text{C}$

$\theta_3 = 10^\circ\text{C}$

 \vdots

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

$\theta_{365} = 2^\circ\text{C}$

Initialize $V_0 = 0$

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

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

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

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

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

 $\beta = 0.9$, V_t as ~~moving~~ averaging over $\approx \frac{1}{1-\beta}$ days' temperature

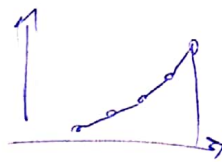
$\frac{1}{1-0.9} = 10$

 $\beta = 0.9$: average over 10 days

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

 $\beta = 0.98$: average over 50 days $\beta = 0.5$ average over 2 days

\rightarrow plot is smoother with higher β but curve is actually shifted (latency)
 \rightarrow much more noisy but adapts quicker to temperature changes (less latency)

 $\Rightarrow \beta$ is also hyperparameter

$$\frac{0.9}{1-\beta} \approx 0.35 \approx \frac{1}{e} \quad \frac{(1-\epsilon)^{1/\epsilon}}{0.9} = \frac{1}{e}$$

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

Implementation

$V_0 = 0$

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

$V_1 := \beta \cdot V_0 + (1-\beta) \theta_2$

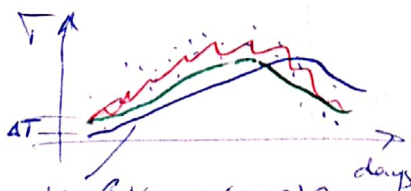
 \vdots

repeat {

get next θ_t

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

}

Bias correction

$V_t = \beta V_{t-1} + (1-\beta) \theta_t$
 \hookrightarrow Starts with lower value at the beginning!

 \hookrightarrow do this:

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

 \hookrightarrow get better estimation for early values

Gradient Descent with momentum → compute exponentially ^{weighted} averaged of the gradients and use that to update the weights

GD with momentum



oscillations slow down GD and prevent us from taking a large learning rate

↑ want slower learning on vertical axis
→ want faster learning on horizontal axis

GD with momentum will take smaller steps in vertical direction and converge faster → damping out oscillations

with momentum:

On iteration t : Initialize $V_{dw}=0, V_{db}=0$

Compute dW, db on current mini-batch (works also for BGD)

$$\begin{aligned} V_{dw} &= \beta V_{dw} + (1-\beta) dW \\ V_{db} &= \beta V_{db} + (1-\beta) db \end{aligned}$$

(moving average of gradient) derivative
like friction ↑ like a velocity acceleration

Intuition:

Analogy: ball rolling down a hill

update parameters:

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

→ smoothing out steps of gradient descent

PYTHON

np.zeros_like()

Implementation:

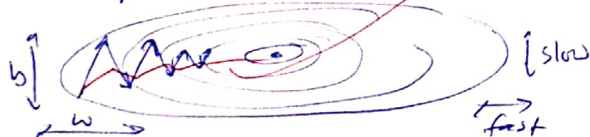
Hyperparameters: α, β

most common $\beta = 0.9$ = averaging over last 10 gradients

In practice: don't really see bias correction when implementing GD with momentum

RMSprop (Root mean square prop)

example from above:



with RMSprop (could use larger learning rate α)

On iteration t :

Compute dW, db on current mini-batch

$$\begin{aligned} S_{dw} &= \beta_2 S_{dw} + (1-\beta_2) dW^2 \\ S_{db} &= \beta_2 S_{db} + (1-\beta_2) db^2 \end{aligned}$$

element-wise
keeping exponentially weighted average of squares of derivatives

update parameters

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

→ derivatives are larger in "vertical" than in "horizontal" direction

add ϵ to denominator $\approx 10^{-8}$
→ for numerical stability

divide by small number helps damping out the oscillations

divide by large number

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

ADAM optimization algorithm \rightarrow momentum + RMSprop

\hookrightarrow works well as generalized opt. algos.

β_1 & $\beta_2 \rightarrow$ hyperparameters

initialize:

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

On iteration t :

compute dw, db using current mini-batch (MBGD)

momentum exponentially averaged \nearrow

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

\rightarrow momentum-like update with β_1

RMSprop \rightarrow

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

\rightarrow RMSprop-like update with β_2

Implementation:

\rightarrow after bias correction

$$\overset{\text{corrected}}{V_{dw}} = V_{dw} / (1 - \beta_1^t), \quad \overset{\text{corrected}}{V_{db}} = V_{db} / (1 - \beta_1^t)$$
$$\overset{\text{corrected}}{S_{dw}} = S_{dw} / (1 - \beta_2^t), \quad \overset{\text{corrected}}{S_{db}} = S_{db} / (1 - \beta_2^t)$$

Perform update:

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

Hyperparameters choice:

α : needs to be tuned \leftarrow try range

$\beta_1: 0.9$ (dw)

$\beta_2: 0.999$ (dw^2 , db^2)

$\epsilon: 10^{-8}$

Quiz

[3] (7) {2}

a

$\theta_1 = 10^\circ\text{C}$
 $\theta_2 = 10^\circ\text{C}$
 $\beta = 0.5$
 $V_0 = 0$

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

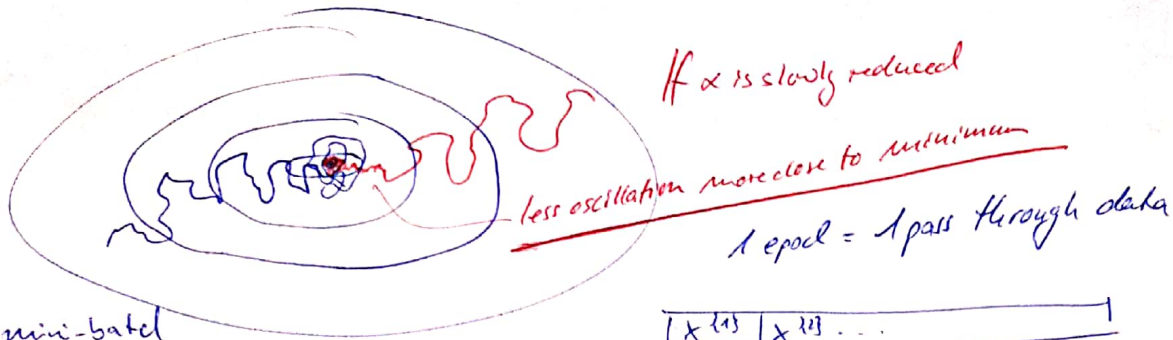
$$V_2 = 0.5 \cdot 10^\circ\text{C} + (0.5) \cdot 10^\circ\text{C} = 10^\circ\text{C}$$

5 + 5

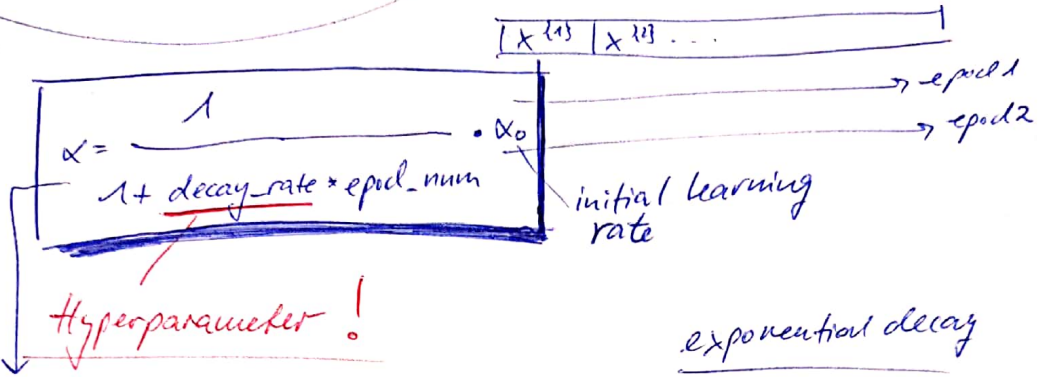
$$V_2^{\text{corr}} = \frac{10^\circ\text{C}}{1 - 0.5^2} = \frac{0.5 \cdot 10^\circ\text{C} + 10^\circ\text{C}}{1 - 0.5^2}$$

(5)

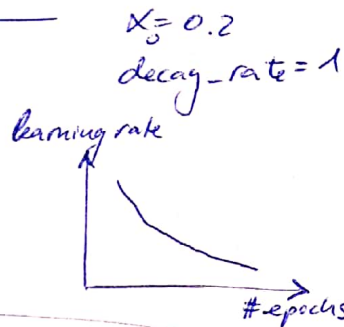
Learning rate decay ^{way to} \rightarrow speed up learning algorithm, is to slowly reduce α over time



mini-batch with fixed α



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



exponential decay

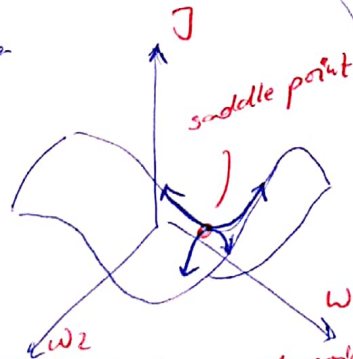
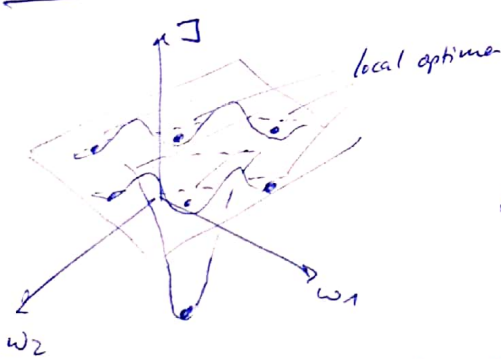
$$\alpha = 0.95^{\text{epoch_num}} \cdot \alpha_0$$

other

$$\alpha = \frac{k}{\sqrt{\text{epoch_num}}} \cdot \alpha_0$$

$$\text{or } \frac{k}{\sqrt{t}} \cdot \alpha_0$$

The problem of local optima



In NN training most points of zero-gradients are not local optima, instead are saddle points!

If gradient = 0 in real direction can either be convex or concave
In 20000 dimensional space ^{more} likely to be in saddle point than in local optimum

Plateau can really slow down learning

regime where gradient too for a long time derivative

- Unlikely to get stuck in local optima as long as NN is big enough, J defined over high dimensional space \rightarrow JER large
- Plateau can make learning slow \rightarrow Adam \rightarrow step size \rightarrow momentum



Assignment Optimization Methods

Gradient Descent

- Difference between batch gradient descent, mini-batch GD and stochastic GD is the number of examples you use to perform one update step
- You have to tune a learning rate hyperparameter α
- A well-tuned mini-batch size usually outperforms GD or SGD, especially when training set is large

Mini-Batch Gradient Descent

- Shuffling and partitioning are the two steps required to build mini-batches
- Powers of 2 are commonly chosen $\rightarrow 16, 32, 64, 128, \dots$

Momentum

- takes past gradients into account to smooth out the steps of GD
It can be applied with BGD, MBGD or SGD
- You have to tune a momentum hyperparameter β and α

Adam