



Optimization Algorithms

Mini-batch gradient descent

Speed up model training.

DL is a highly empirical subject
⇒ skills to speed up training

Batch vs. mini-batch gradient descent

Vectorization allows you to efficiently compute on m examples.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \quad \begin{matrix} (n_x, m) \\ \underbrace{\hspace{10em}}_{X^{[1:3]} \quad (n_x, 1000)} \quad \underbrace{\hspace{10em}}_{X^{[2:3]} \quad (n_x, 1000)} \quad \dots \quad \underbrace{\hspace{10em}}_{X^{[5000:3]} \quad (n_x, 1000)} \end{matrix}$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix} \quad \begin{matrix} (1, m) \\ \underbrace{\hspace{10em}}_{Y^{[1:3]} \quad (1, 1000)} \quad \underbrace{\hspace{10em}}_{Y^{[2:3]} \quad (1, 1000)} \quad \dots \quad \underbrace{\hspace{10em}}_{Y^{[5000:3]} \quad (1, 1000)} \end{matrix}$$

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

5,000 mini-batches of 1,000 each
Mini-batch t : $X^{[t:3]}, Y^{[t:3]}$

$x^{(i)}$ is training example i
 $\approx [x]$ layer $[]$
 $X^{[t:3]}, Y^{[t:3]}$ — batches
— Andrew Ng

$$X = [x^{(1)} x^{(2)} x^{(3)} \dots x^{(m)}]$$

if m is large, it can be slow.
need to process the whole training set
before gradient descent

5000 mini batches, 1000 each

(n_x , size of each batch)

of batches

Mini-batch gradient descent

Repeat for $t = 1, \dots, 5000$ {

Forward prop on $X^{(t)}$.

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

$$A^{(t)} = \sigma^{(t)}(Z^{(t)})$$

} Vectorized implementation (1000 examples)

↙ for $X^{(t)}, y^{(t)}$

Compute cost $J^{(t)} = \frac{1}{1000} \sum_{i=1}^{1000} \ell(w, y_i) + \frac{\lambda}{2 \cdot 1000} \sum_{i=1}^{1000} \|w^{(t)}\|_2^2$ → regularization

Backprop to compute gradients w.r.t $J^{(t)}$ (using $(X^{(t)}, y^{(t)})$)

$W^{(t+1)} = W^{(t)} - \alpha dW^{(t)}, b^{(t+1)} = b^{(t)} - \alpha db^{(t)}$

one pass through training set
"one epoch"

"1 epoch" pass through training set.

Andrew Ng

vectorization to process all 1000 examples in a certain batch

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

$$A^{(t)} = \sigma^{(t)}(Z^{(t)})$$

$$A^{(t)} = \sigma^{(t)}(Z^{(t)})$$

runs much faster. Everyone will use with large datasets

mini-batch2



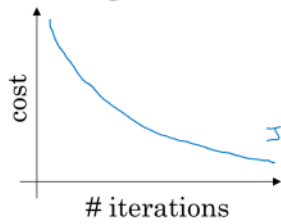
deeplearning.ai

Optimization Algorithms

Understanding mini-batch gradient descent

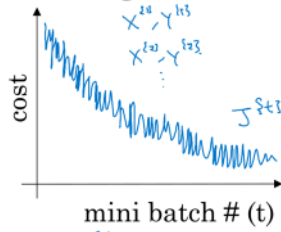
Training with mini batch gradient descent

Batch gradient descent



Always decreases (one every iteration).

Mini-batch gradient descent



Plot J^{t+1} computed using X^{t+1}, Y^{t+1}

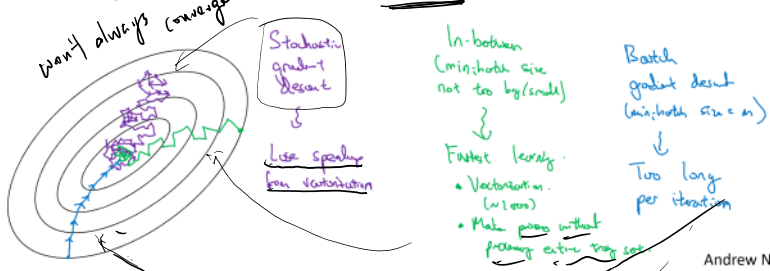
Andrew Ng

Need not always decrease
should trend down.
but noisier!

Choosing your mini-batch size

→ If mini-batch size = m : Batch gradient descent. $(X^{(1)}, Y^{(1)}) = (X, Y)$.
→ If mini-batch size = 1: Stochastic gradient descent. Every example is its own mini-batch. $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$ mini-batch.

In practice: Somewhere in-between 1 and m



Andrew Ng

two extremes. 每 batch → 全部数据
→ 一条数据

choose something in between; not 1, not m .

Choosing your mini-batch size

If small toy set: Use batch gradient descent. ($m \leq 2000$)

Typical mini-batch sizes:

→ 64, 128, 256, 512

$\frac{1024}{2^{10}}$

Make sure mini-batch fit in CPU/GPU memory. $X^{(1)}, Y^{(1)}$

REMEMBER THESE RULES.

$m \leq 2000$ batch.

7200 sizes of mini-batch $2^6 = 64$
因为内存 $2^7 = 128$ 最常见.
size 更快 $2^8 = 256$
(不浪费空间) $2^9 = 512$
 $2^{10} = 1024$ 有更大.

Don't exceed CPU/GPU memory.

Andrew Ng

Approach: Try different powers of 2. Use the most efficient one (size).

1. ... 1. mini-batch size and

up next: none efficient than mini-batch gd. ^{the most efficient one (size)}

exp weighted average



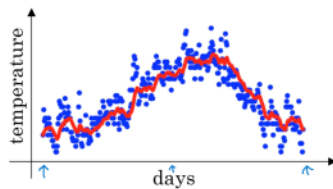
deeplearning.ai

Optimization Algorithms

Exponentially weighted averages

Temperature in London

$\theta_1 = 40^\circ\text{F}$ $4^\circ\text{C} \leftarrow$
 $\theta_2 = 49^\circ\text{F}$ 9°C
 $\theta_3 = 45^\circ\text{F}$
 \vdots
 $\theta_{180} = 60^\circ\text{F}$ 15°C
 $\theta_{181} = 56^\circ\text{F}$
 \vdots

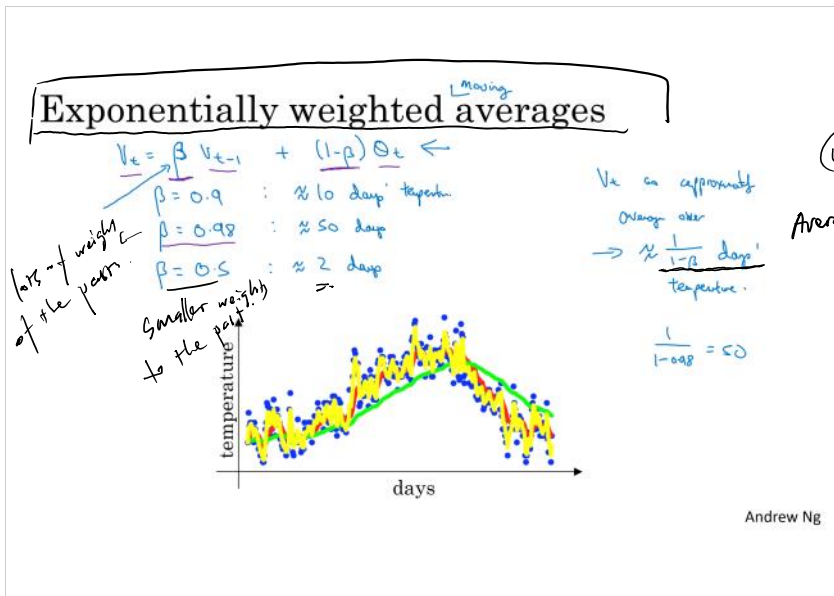


$$\begin{aligned}
 v_0 &= 0 \\
 v_1 &= 0.9 v_0 + 0.1 \theta_1 \\
 v_2 &= 0.9 v_1 + 0.1 \theta_2 \\
 v_3 &= 0.9 v_2 + 0.1 \theta_3 \\
 &\vdots \\
 v_t &= 0.9 v_{t-1} + 0.1 \theta_t
 \end{aligned}$$

AKMA(1, 1) ?

non-stationary s.f.

Andrew Ng



why? 能算出来吧?
 Averaging over the last $\frac{1}{1-\beta}$ days temperature.

Next what this is doing

ewa2



Optimization Algorithms

Understanding exponentially weighted averages

red. green. — yellow.

Test a rule of thumb of how to think about it not mathematical derivation!!

to test
add up to 1
bias correction

$$\zeta = 0.1 \quad (1 - 0.02)^{0.02} = 0.98^{0.02} \approx \frac{1}{2} \text{ about } \frac{1}{2}$$

(?)

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

怎么计算结果的呢?

这部分是用来算 Average 的

$$\begin{aligned} n^3 &= \beta n^5 + (1 - \beta) \theta^3 \\ n^5 &= \beta n^7 + (1 - \beta) \theta^5 \\ n^7 &= \beta n^9 + (1 - \beta) \theta^7 \\ n^9 &= 0 \end{aligned}$$

averages

Implementing exponentially weighted

$$\begin{aligned} n^0 &:= \beta n^0 + (1 - \beta) \theta^0 \leftarrow \text{repeat} \\ &\vdots \\ n^0 &:= \beta n^0 + (1 - \beta) \theta^0 \\ n^0 &:= 0 \end{aligned}$$

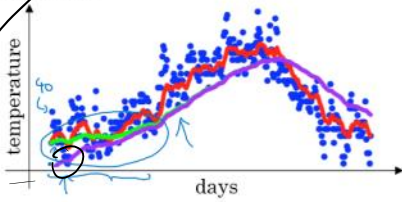
Advantage: takes very little memory efficient. takes one line of code. Not the most accurate way to compute the average



Optimization Algorithms

Bias correction in exponentially weighted average

Bias correction



$$\begin{aligned} \rightarrow v_t &= \beta v_{t-1} + (1 - \beta) \theta_t \\ v_0 &= 0 \\ v_1 &= 0.98 v_0 + 0.02 \theta_1 \\ v_2 &= 0.98 v_1 + 0.02 \theta_2 \\ &= 0.98 \times 0.02 \theta_1 + 0.02 \theta_2 \\ &= 0.0196 \theta_1 + 0.02 \theta_2 \end{aligned}$$

$$\begin{aligned} \frac{v_t}{1 - \beta^t} \\ t=2: 1 - \beta^t &= 1 - (0.98)^2 = 0.0396 \\ \frac{v_2}{0.0396} &= \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396} \end{aligned}$$

Andrew Ng

β large \Rightarrow the initial value $v_0 = 0$ keeps the predicted value v low for a while, not good..

to modify this to make it better

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

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

As $t \uparrow$, $\Rightarrow \beta^t \Rightarrow 0$ this correction diminishes

only helpful for the "warm-up" period.

But 0 people don't really bother to do this in ML



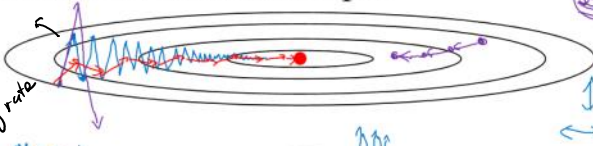
Optimization Algorithms

Gradient descent with momentum

works faster than gradient descent.

Gradient descent example

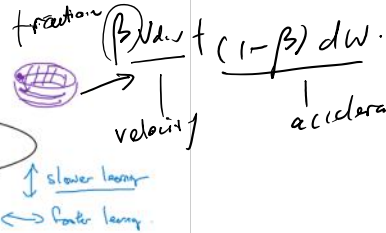
slow down learning
can't use larger learning rate



Momentum:
On iteration t :
Compute dW, db on current mini-batch.
 $v_{dw} = \beta v_{dw} + (1-\beta) dW$
 $v_{db} = \beta v_{db} + (1-\beta) db$
Direction \rightarrow velocity
 $W = W - \alpha v_{dw}$
deelearning.ai



$$v_0 = \beta v_0 + (1-\beta) dW$$



"ball rolling down a bowl analogy"

Moving average of the derivatives of W .

What does average out (cancel out)

in vertical $\frac{d^2}{dt^2}$ cancel out
horizontal $\frac{d^2}{dt^2}$ won't cancel out.

Andrew Ng

Want: larger learning rate horizontally.

Implementation details

$$v_{dw} = 0, v_{db} = 0$$

On iteration t :

Compute dW, db on the current mini-batch

$$\begin{aligned} \rightarrow v_{dw} &= \beta v_{dw} + (1-\beta) dW \\ \rightarrow v_{db} &= \beta v_{db} + (1-\beta) db \\ W &= W - \alpha v_{dw}, \quad b = b - \alpha v_{db} \end{aligned}$$

in the literature. A different type of parametrization.
Differ? Scale α just affect α 's choice.
this way of parametrization is less intuitive.

Hyperparameters: α, β

$\beta = 0.9$ → appear robust.
average over last 10 gradients

Andrew Ng

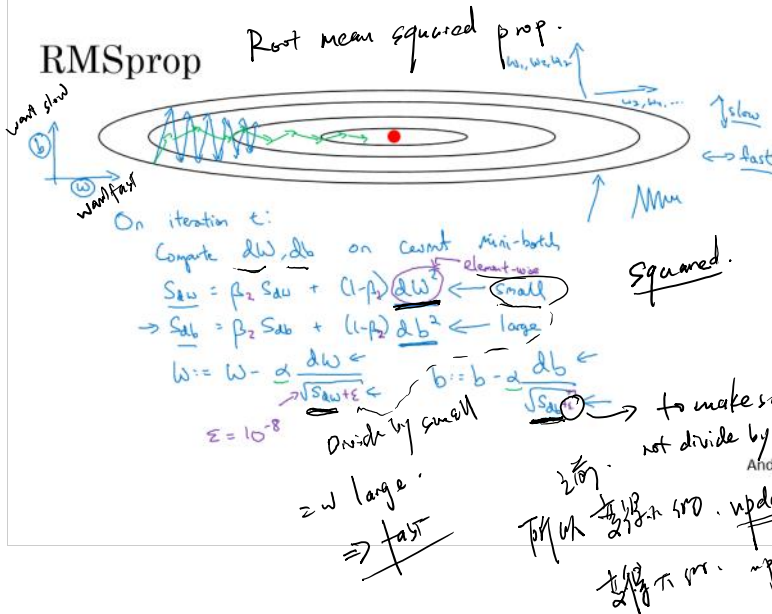
$$\beta = 0.9 \quad 0.9^{\frac{1}{0.1}} = 0.9^{10} = \frac{1}{e}$$

bias correction? Don't bother.
after just 10 iteration. the moving average would have been warmed up.
the initial value won't cause much trouble.



Optimization Algorithms

RMSprop



exponentially weighted average of squared derivatives.

In practice, high dimension

this is just a 2-dim example

Next video : Combine Momentum & RMSprop. \Rightarrow Adam optimization



Optimization Algorithms

Adam optimization algorithm

RMSprop & Adam are the two optim algorithms that stand out among many new optimizer algorithms

Proven effective

Adam optimization algorithm

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

On iter t :

Compute dw, db using current mini-batch

$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw, V_{db} = \beta_1 V_{db} + (1-\beta_1) db \leftarrow \text{"moment"} \beta_1$

$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2, S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \leftarrow \text{"RMSprop"} \beta_2$

yhat = np.array([.9, 0.2, 0.1, .4, .9])

bias correction

$$\begin{cases} \frac{V_{dw}^{corrected}}{S_{dw}^{corrected}} = \frac{V_{dw}}{(1-\beta_1^t)}, & \frac{V_{db}^{corrected}}{S_{db}^{corrected}} = \frac{V_{db}}{(1-\beta_1^t)} \\ \frac{S_{dw}^{corrected}}{S_{db}^{corrected}} = \frac{S_{dw}}{(1-\beta_2^t)}, & \frac{S_{db}^{corrected}}{S_{db}^{corrected}} = \frac{S_{db}}{(1-\beta_2^t)} \end{cases}$$

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

Andrew Ng

- 1/10 rate ...

Adaptive moment estimate.

Hyperparameters choice:

- $\rightarrow \alpha$: needs to be tune
- $\rightarrow \beta_1$: 0.9 $\rightarrow (dw)$
- $\rightarrow \beta_2$: 0.999 $\rightarrow (dw^2)$
- $\rightarrow \epsilon$: 10^{-8} \leftarrow Doesn't matter

\leftarrow try a range of α to see which works

much, no need to tune

1st is 1st second moment is β ?

Adam: Adaptive moment estimation



Adam Coates

Andrew Ng

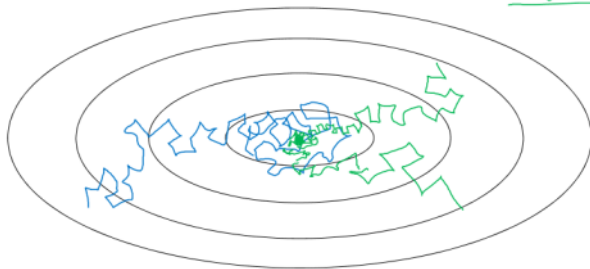
not related to this guy.



Optimization Algorithms

Learning rate decay

Learning rate decay



Andrew Ng

Learning rate decay

1 epoch = 1 pass through data.

$$\alpha' = \frac{1}{1 + \text{decay_rate} * \text{epoch_num}} \alpha_0$$

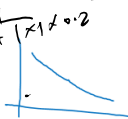
Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
\vdots	\vdots

$$\alpha_2 = \frac{1}{1+1} \times 0.2$$



$$\alpha_0 = 0.2$$

$$\text{decay_rate} = 1$$



Andrew Ng

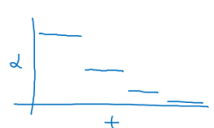
slowly reducing learning rate
start with high LR.
then when it is close to
convergence \Rightarrow reduce LR.

Other learning rate decay methods

看起来都好 ad hoc...

Formula

$$\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \quad \text{-- exponentially decay.}$$

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


discrete staircase

interesting. why not continuous?

Manual decay if training only one model.

works if training ^{Andrew Ng} only a few number of model.

Next week: how to systematically tune param.

LR decay is not the top priority for param tuning

不是调 LR 的 top priority AN.

local optimal

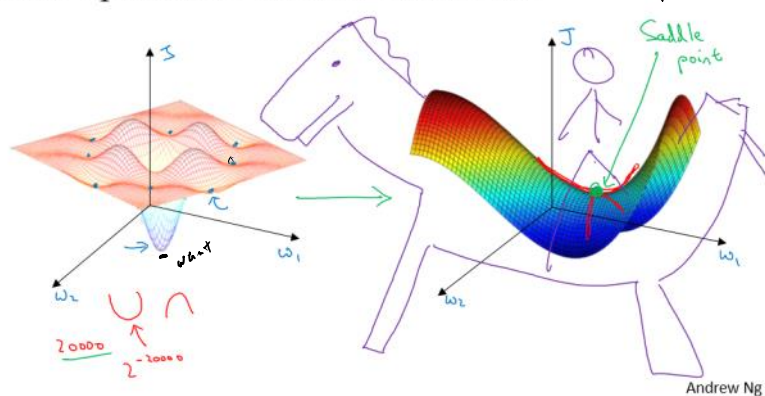


deeplearning.ai

Optimization Algorithms

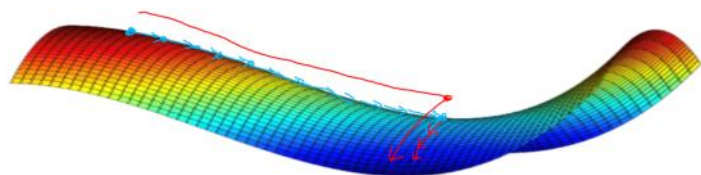
The problem of local optima

Local optima in neural networks



most points of zero gradient are not optimal.
They are saddle points.
• saddle points problematic for high dim data

Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow

plateau
slow down training
take a long time to get to optimal.

Take away:

- local optima not much a problem for DL tasks
- plateaus slow down learning

Andrew Ng

$$V_1 = 0.5 \times 10 + 0.5 \times 10 = 5$$

$$V_2 = 0.5 \times 15 + 0.5 \times 10 = 7.5$$

• key: high dim spaces.
understanding on high-dim space
is still evolving.