

Week 2: Practical Aspects of Deep Learning

Mini Batch Gradient Descent

- Let's start looking at mini batch gradient descent by looking at the X and Y matrices. X contains all the training examples and is therefore equal to $[x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}]$. Y is also similar but it contains the labels so it is equal to $[y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(m)}]$. We live in the age of big data so it is common to have large datasets; imagine if our data set has 5 million training examples. Batch gradient descent will need to use all 5 million examples on every step making this process very slow. mini-batch gradient descent, on the other hand, only requires a few examples 'per batch'; in the example below we are using batches of 1,000 examples so for X we have $X^{(1)} = \text{the first 1,000 examples}$ and it has dimensions $(n_x, 1000)$, $X^{(2)} = [x^{(1001)}, \dots, x^{(2000)}]$ and so on. Y goes through the same process. In our example we then end up with 5,000 mini batches of 1,000 examples each. Each mini batch t then is composed of $X^{(t)}, Y^{(t)}$.

Batch vs. mini-batch gradient descent

Vectorization allows you to efficiently compute on m examples.

$$\begin{aligned}
 X &= \underbrace{[x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(1000)}]}_{X^{(1)} \ (n_x, 1000)} \underbrace{[x^{(1001)} \ \dots \ x^{(2000)}]}_{X^{(2)} \ (n_x, 1000)} \dots \underbrace{[\dots \ x^{(m)}]}_{X^{(5,000)} \ (n_x, 1000)} \\
 Y &= \underbrace{[y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)}]}_{Y^{(1)} \ (1, 1000)} \underbrace{[y^{(1001)} \ \dots \ y^{(2000)}]}_{Y^{(2)} \ (1, 1000)} \dots \underbrace{[\dots \ y^{(m)}]}_{Y^{(5,000)} \ (1, 1000)}
 \end{aligned}$$

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

5,000 mini-batches of 1,000 each

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

$x^{(1)}$
 $z^{[2]}$
 $X^{(t)}, Y^{(t)}$

Andrew Ng

- How do we implement mini batch gradient descent? We start with a for-loop for all the training batches, using our previous examples we had 5,000 so that's the index for the for-loop, we then compute Z as $Z^{(l)} = w^{(l)}X^{(t)} + b^{(l)}$, A, until we get $A^{[L]}$. We then compute the cost, do backprop to compute our gradients using $(X^{(t)}, Y^{(t)})$ and update our gradients as follows: $w^{[l]} = w^{[l]} - \alpha dw^{[l]}$ and $b^{[l]} = b^{[l]} - \alpha db^{[l]}$. This is one pass of gradient descent, also known as epoch.

- Summarizing, mini batch gradient descent is the preferred choice when we have a large training set. It runs much faster than gradient descent and is what everybody uses when training on a large data set.

Mini-batch gradient descent

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

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

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

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

$$\vdots$$

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

Compute cost $J^{(t)} = \frac{1}{1000} \sum_{i=1}^{1000} d(\hat{y}^{(t)}, y^{(t)}) + \frac{\lambda}{2 \cdot 1000} \sum_{i=1}^{1000} \|W^{(t)}\|_F^2$

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

$$W^{(t+1)} = W^{(t)} - \alpha \frac{\partial J}{\partial W}, \quad b^{(t+1)} = b^{(t)} - \alpha \frac{\partial J}{\partial b}$$

}

"1 epoch"
 — pass through training set.

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

X, Y

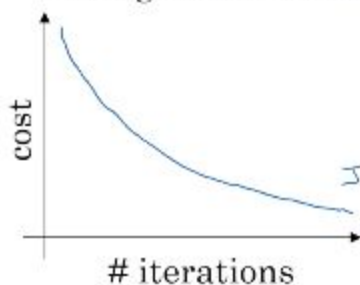
Andrew Ng

Understanding Mini Batch Gradient Descent

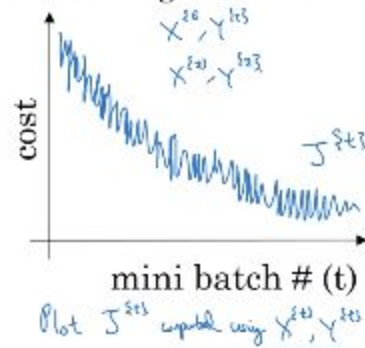
- When we train with regular batch gradient descent we expect to see a cost reduction every time. If we see an increase even in one pass we know there is something wrong. That is not the case with mini-batch gradient descent. Mini batch gradient descent is noisier because it is like training on a different training set every time. So one batch might be an 'easy' batch $X^{(t)}, Y^{(t)}$ and another batch $X^{(t)}, Y^{(t)}$ might be 'harder' (meaning it might contain some unlabeled examples). In general mini batch gradient descent should also trend lower over time but it will be noisier.

Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent



Andrew Ng

- How do we choose the mini batch size? Let's start by looking at the extremes: if we have minibatch size of m then we have batch gradient descent where $(X^{(t)}, Y^{(t)}) = (X, Y)$. If we have minibatch size = 1 then we have stochastic gradient descent where each training example is its own mini batch. In practice we choose something in between.

Choosing your mini-batch size

→ If mini-batch size = m : Batch gradient descent. $(X^{(t)}, Y^{(t)}) = (X, Y)$
 → If mini-batch size = 1 : Stochastic gradient descent. Every example is its own mini-batch. $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(n)}, y^{(n)})$ mini-batch.
 In practice: Somewhere in-between 1 and m



In general, batch gradient descent takes a straightforward to the global minimum, which is good, but it takes too long per iteration because each iteration uses the whole training set. Stochastic gradient descent takes a very noisy road to the global minimum and will never converge, rather it will 'hang out' around that global minimum. The disadvantage of stochastic gradient descent is that we lose the speed from vectorization. Choosing a batch size of somewhere in between enables the fastest learning because we are still using vectorization and we don't need to wait to process the entire training set.

- Here are some guidelines for how to choose the mini batch size:
 - If we have a small training set: use batch gradient descent
 - Typical batch sizes are multiples of 2 (64, 128, 256, 512) because of the way computer memory is laid out and accessed
 - Finally, make sure the mini batch size fits in CPU/GPU memory

Choosing your mini-batch size

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

Typical mini-batch sizes:

64, 128, 256, 512
 $2^6, 2^7, 2^8, 2^9$

1024
 2^{10}

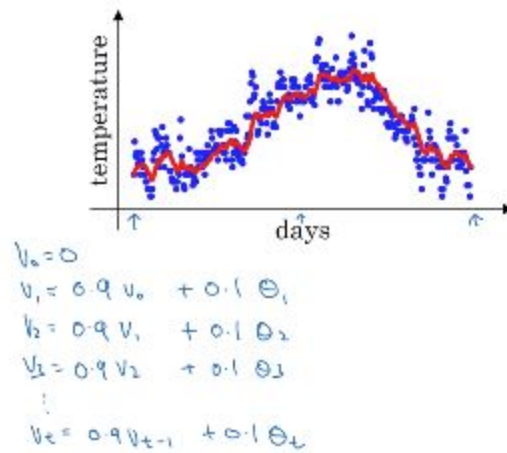
Make sure mini-batch size fits in CPU/GPU memory.
 $x^{(1)}, y^{(1)}$

Exponentially Weighted Averages

- Now we will take a look at other optimization algorithms that are faster than gradient descent. In order to understand these algorithms we need to understand something called exponentially weighted averages. We'll illustrate it using the example below. We have the temp in London for one year. When we plot it we have something like the blue dots. We then do the following: set $v_0 = 0$, $v_1 = 0.9 * v_0 + 0.1 * \theta_1$ where v_0 is the previous day v and θ_1 is the temperature that day. The general formula is then $v_t = 0.9 * v_{t-1} + 0.1 * \theta_t$. If we plot that we then end up with the red line below, a less noisier line. This is what is called an exponentially weighted average of the daily temperature.

Temperature in London

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



Andrew Ng

- The general formula is then $v_t = \beta * v_{t-1} + (1 - \beta) * \theta_t$. We can also think of v_t as approximately averaging over $1/(1 - \beta)$ days' temperature. If $\beta = 0.9$ we are averaging over ≈ 10 days' temp. If $\beta = 0.98$ we are averaging over ≈ 50 days' temp and our plot is the green line below. If $\beta = 0.5$ we are averaging over ≈ 2 days' temp and our plot is the yellow one below,

Exponentially weighted averages ^{↳ Moving}

$$v_t = \beta v_{t-1} + (1-\beta) \theta_t \leftarrow$$

$\beta = 0.9$: ≈ 10 days' temp

$\beta = 0.98$: ≈ 50 days

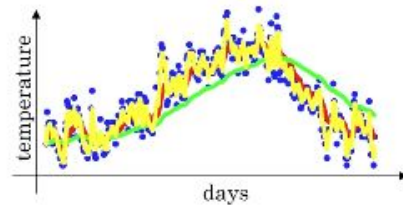
$\beta = 0.99$: ≈ 100 days

v_t is approximately

average of

$\rightarrow \approx \frac{1}{1-\beta}$ days' temperature.

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



Averaging over a small number of days adapts fast to a new day temp but is very noisy. Averaging over a big number of days gives us a smoother curve but gives less weight to a new day's temp.

Understanding Exponentially Weighted Averages

- Let's go a little deeper into how exponentially weighted averages work. Remember our equation is $v_t = \beta v_{t-1} + (1-\beta) \theta_t$. If we take that equation and plug in the values from day 100 we see the equations on the left side below. We can then do some algebra using $\beta = 0.9$, we know $v_{100} = 0.1 \theta_{100} + 0.9 v_{99}$ but we can substitute v_{99} with its equation so we then have $v_{100} = 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 v_{98})$. We can then substitute v_{98} with its equation and so on. This equation ends up as a weighted sum/average of the day's temp + and the previous day temp and the previous day temp and so on.

Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$\begin{aligned}
 v_{100} &= 0.9v_{99} + 0.1\theta_{100} \\
 v_{99} &= 0.9v_{98} + 0.1\theta_{99} \\
 v_{98} &= 0.9v_{97} + 0.1\theta_{98} \\
 &\dots \\
 \rightarrow v_{100} &= 0.1\theta_{100} + 0.9(0.1\theta_{99} + 0.9(0.1\theta_{98} + 0.9(0.1\theta_{97} + \dots))) \\
 &= 0.1\theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1(0.9)^2 \theta_{98} + 0.1(0.9)^3 \theta_{97} + \dots \\
 &\quad \underbrace{0.9^{10} \approx 0.35 \approx \frac{1}{e}} \quad \underbrace{(1-\epsilon)^{1/\epsilon} = \frac{1}{e}}_{\epsilon=0.01 \rightarrow 0.98^{50} \approx \frac{1}{e}}
 \end{aligned}$$

Andrew Ng

We can draw this equation on the top right corner and we end up with an exponentially decaying function. How many days is this equation averaging over? It turns out that $0.9^{10} \approx 0.35 \approx 1/e$. This means it takes about 10 days for the height of this to decay $\frac{1}{e}$ so we say we are computing the average that focuses on the last 10 days because after 10 days the weight decays to less than about a $\frac{1}{e}$ of the weight of the current day.

- How do we actually implement exponentially weighted averages? What we do in practice is shown on the right side below. We initialize $v_0 = 0$ then update it to $v_0 = \beta v + (1 - \beta) * \theta_1$ and so on. We can do the same using a for loop.

Implementing exponentially weighted averages

$$\begin{aligned}
 v_0 &= 0 \\
 v_1 &= \beta v_0 + (1 - \beta) \theta_1 \\
 v_2 &= \beta v_1 + (1 - \beta) \theta_2 \\
 v_3 &= \beta v_2 + (1 - \beta) \theta_3 \\
 &\dots
 \end{aligned}$$

$$\begin{aligned}
 v_0 &= 0 \\
 v_0 &:= \beta v + (1 - \beta) \theta_1 \\
 v_0 &:= \beta v + (1 - \beta) \theta_2 \\
 &\vdots
 \end{aligned}$$

```

→ v0 = 0
Repeat {
  Get next θt
  v0 := β v0 + (1-β) θt ←
}
    
```

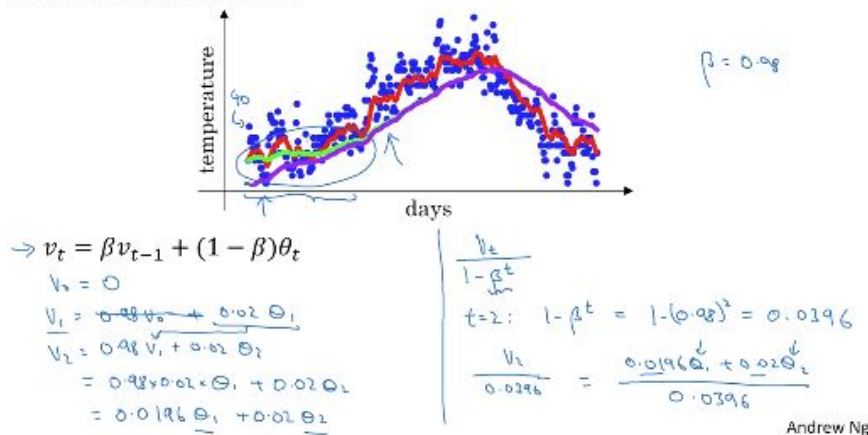
Andrew

Exponentially weighted averages is one line of code and very efficient from a computation and memory point of view which is one reason why it is commonly used in ML.

Bias Correction in Exponentially Weighted Averages

- In the last session we saw how different beta values give us different curves as shown below. It turns out we do not obtain the green line below, we get the purple line. Notice how this line is defaced to the right. The reason for this is shown on the bottom left corner: the computation for day 1, θ_1 , is a low value and not a good average; we get the same case for day 2. We correct this with bias correction. What we do is the following: $v_t/(1 - \beta^t)$. If $t=2$ then we have $1 - \beta^t = (1 - 0.98)^2 = 0.0396$ so $v_t/0.0396$ will give us a better average.

Bias correction



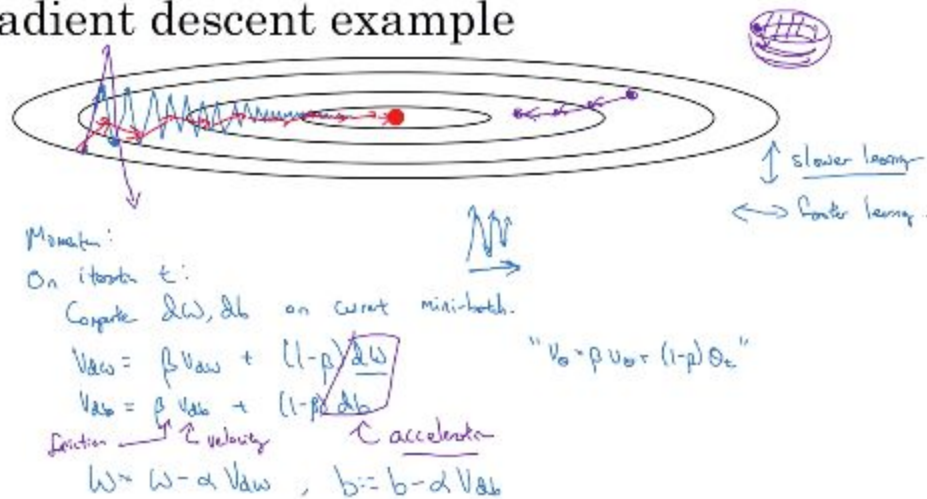
Gradient Descent with Momentum

- We have an algorithm called gradient descent with momentum that works faster than normal gradient descent. The basic idea is to compute an exponentially weighted average of the gradients and use that instead. Let's see how it works: gradient descent usually has a somewhat erratic way to get to the global minimum as shown below. What we want with gradient descent is slower learning vertically and faster learning horizontally; in other words, a smoother and faster path to the global minimum.

Momentum provides this in the following manner:

- On each iteration t , we compute dw and db as always
- We then compute $v_{dw} = \beta * v_{dw} + (1 - \beta)dw$ and $v_{db} = \beta * v_{db} + (1 - \beta)db$
- We then update our weights using v_{dw} and v_{db} : $w = w - \alpha * v_{dw}$, $b = b - \alpha v_{db}$

Gradient descent example



- Let's look at a couple of implementation details: below are the steps needed to do gradient descent with momentum. In practice $\beta = 0.9$ is a good value to use and it means that we are averaging over the last 10 gradients. Notice we now have an additional hyperparameter, β , which controls the exponentially weighted average.

Implementation details

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

On iteration t :

Compute dW, db on the current mini-batch

$$\begin{aligned} \rightarrow v_{dw} &= \beta v_{dw} + (1-\beta) dW \\ \rightarrow v_{db} &= \beta v_{db} + (1-\beta) db \end{aligned} \quad \left| \quad v_{dw} = \beta v_{dw} + dW \leftarrow$$

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

$$\frac{v_{dw}}{\beta^t}$$

Hyperparameters: α, β

$$\beta = 0.9$$

average over last 10 gradients

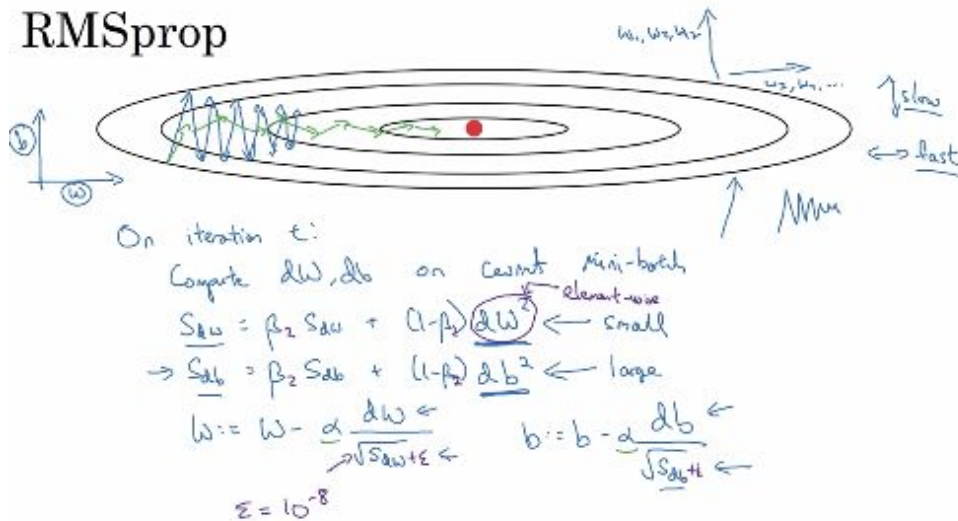
Ar

Often we will see v_{dw} as $v_{dw} = \beta v_{dw} + dW$. This also works fine but Professor Ng prefers the other version because this one is less intuitive. In practice we also see people not bothering with bias correction.

RMS Prop

- There is another algorithm that can speed up gradient descent: Root Mean Square Prop or RMSProp. Similar to Momentum, RMS prop smooths out the gradient descent steps by trying to minimize the vertical learning and maximizing the horizontal learning. The steps are:
 - On each iteration t , compute dw and db as always

- We then compute $s_{dw} = \beta * s_{dw} + (1 - \beta)dw^2$ and $s_{db} = \beta * s_{db} + (1 - \beta)db^2$
- We then update our gradients as follows:
 $w = w - \alpha * (dw / \sqrt{s_{dw}})$ and $b = b - \alpha * (db / \sqrt{s_{db}})$



Andrei

The idea is that dw is a small value and db is a large value so when we update the gradients using those values we get the smooth steps we are looking for.

Adam Optimization Algorithm

- Adam Optimization basically combines Momentum & RMS Prop. The calculations are:
 - For iteration t we compute dw and db as always
 - We then compute
 $v_{dw} = \beta_1 * v_{dw} + (1 - \beta_1)dw$, $v_{db} = \beta_1 * v_{db} + (1 - \beta_1)db$, $s_{dw} = \beta_2 * s_{dw} + (1 - \beta_2)dw^2$
 and $s_{db} = \beta_2 * s_{db} + (1 - \beta_2)db^2$
 - Here we implement bias correction so we end up with
 $v_{dw}^{corrected} = v_{dw} / (1 - \beta_1^t)$, $v_{db}^{corrected} = v_{db} / (1 - \beta_1^t)$, $s_{dw}^{corrected} = s_{dw} / (1 - \beta_2^t)$, $s_{db}^{corrected} = s_{db} / (1 - \beta_2^t)$
 - We then update our gradients as follows: $w = w - \alpha * (v_{dw}^{corrected} / \sqrt{s_{dw}^{corrected} + \epsilon})$ and
 $b = b - \alpha * (v_{db}^{corrected} / \sqrt{s_{db}^{corrected} + \epsilon})$

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

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

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

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

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

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

Andrew Ng

- This algorithm has a few hyperparameters to use:

- α : needs to be tuned
- β_1 : 0.9
- β_2 : 0.999
- ϵ : 10^{-8}

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}

Adam: Adaptive moment estimation

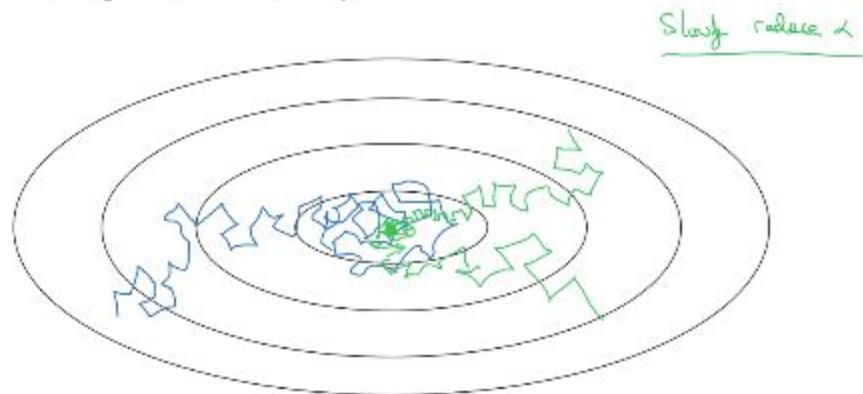
Adam, by the way, stands for Adaptive moment estimation.

Learning Rate Decay

- One of the things that can help our learning algorithm is to slow down the learning rate over time, why? Consider the example below. When doing mini batch gradient descent, we have a noisy way to the global minimum and we never converge. By slowly reducing the learning rate we can have somewhat big steps at the beginning, as shown in the

green trajectory, and smaller ones at the end so our algorithm is tighter and 'hangs out' closer to the global minimum.

Learning rate decay



- How do we implement learning rate decay? Remember that 1 epoch = 1 pass through the training set. What we do is we set $\alpha = 1/(1 + \text{decay rate} * \text{epoch number})$. In this way we get a learning rate that decreases as we go through the epochs, which is exactly what we want.

Learning rate decay

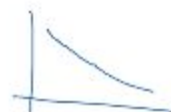
1 epoch = 1 pass through data.

$$\alpha' = \frac{1}{1 + \text{decay-rate} * \text{epoch-number}} \alpha_0$$

Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
...	...



$\alpha_0 = 0.2$
decay rate = 1



- There are other ways of calculating the learning rate decay: other formulas and manually. Other formulas are: $\alpha = 0.95^{\text{epoch number}} * \alpha$, $\alpha = (k/\sqrt{\text{epoch number}}) * \alpha$ or a staircase where we cut alpha in half after a number of steps. Manual decay is done when you have a small number of models that take days to train and you are observing them and changing the learning rate as appropriate (i.e. the learning rate has slowed down, let's cut the value of alpha).

Other learning rate decay methods

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

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

A graph showing α vs t with a staircase-like decay pattern.

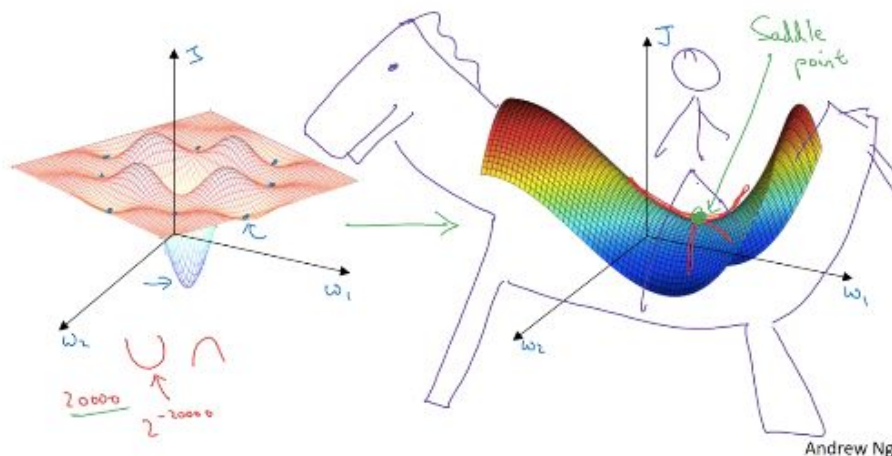
Manual decay.

Ar

The problem of local optima

- In the early days of deep learning there was a lot of worry about local optima and optimization algorithms getting stuck there. As deep learning has advanced this is less of a concern. The picture on the left is what we had in mind when we were afraid of getting stuck in local optima. However, the picture on the right is a better representation of what happens with neural networks. In such a case it is more likely you encounter a 'saddle' point (where the derivative = 0) rather than a local optima. Neural networks usually have many, many parameters so the chances of getting a local optima are very small. The lesson is then that a lot of our intuitions about low dimensional spaces do not translate to high dimensional situations. Therefore, local optima is less of a problem.

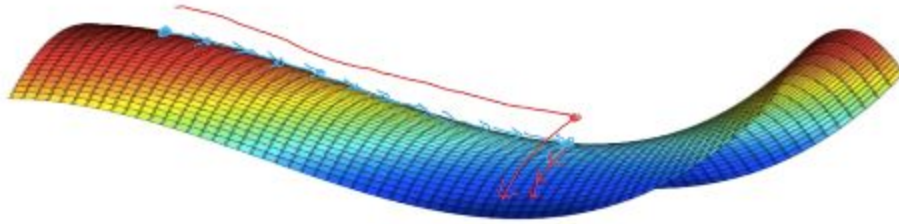
Local optima in neural networks



- If local optima is not a problem, what is a problem? Plateaus. Plateaus can slow down learning and are regions where the derivative is close to 0 for a long time. As shown

below, an optimization algorithm might take a while to get off the plateau and continue to go down towards the global minimum. More sophisticated algorithms such as Adam or RMSprop help us to go quickly through the plateau.

Problem of plateaus



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

Summarizing: local optima is not a problem in deep learning and plateaus can make learning slow.

Summary from programming exercise: Gradient Descent

- The difference between gradient descent, mini batch gradient descent and stochastic gradient descent is the number of examples used in each step.
- You have to tune α .
- A well-tuned mini batch size usually outperforms gradient descent and stochastic gradient descent, especially when the training set is large.

Summary from programming exercise: Mini-Batch Gradient Descent

- Shuffling & Partitioning are two steps required to create mini batches.
- Mini batch sizes are usually powers of two: 64, 128, 256, 512.

Summary from programming exercise: Momentum

- Momentum takes past gradients into account to smooth out the steps of gradient descent. It can be applied to all 3 gradient descent types.
- With momentum you have to tune β and α .

Summary from programming exercise: Momentum

- Some advantages of Adam are relatively memory requirements & that it usually works well even with little tuning of hyperparameters (except α).