



# [2주차] Mini-batch gradient descent

## Batch vs. mini-batch gradient descent

Vectorization allows you to efficiently compute on  $m$  examples.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} & | & x^{(1001)} & \dots & x^{(2000)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix}$$

$(n_x, m)$        $X^{\{1\}} (n_x, 1000)$        $X^{\{2\}} (n_x, 1000)$        $X^{\{5,000\}} (n_x, 1000)$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} & | & y^{(1001)} & \dots & y^{(2000)} & | & \dots & | & \dots & y^{(m)} \end{bmatrix}$$

$(1, m)$        $Y^{\{1\}} (1, 1000)$        $Y^{\{2\}} (1, 1000)$        $Y^{\{5,000\}} (1, 1000)$

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

5,000 mini-batches of 1,000 each

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

$$\begin{matrix} x^{(i)} \\ z^{[l]} \\ X^{\{t\}}, Y^{\{t\}} \end{matrix}$$

Andrew Ng

여기서는 중괄호  $t$ 를 이용해서 mini batches들을 인덱싱 합니다.

# Mini-batch gradient descent

repeat {  
 for  $t = 1, \dots, 5000$  {  
 Forward prop on  $X^{t+1}$ .  
 $Z^{t+1} = W^{t+1} X^{t+1} + b^{t+1}$   
 $A^{t+1} = g^{t+1}(Z^{t+1})$   
 $\vdots$   
 $A^{t+1} = g^{t+1}(Z^{t+1})$  } Vectorial implementation (1000 examples)  
 Compute cost  $J^{t+1} = \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{i=1}^{1000} \|W^{t+1}\|_F^2$ .  
 Backprop to compute gradients w.r.t  $J^{t+1}$  (using  $X^{t+1}, Y^{t+1}$ )  
 $W^{t+1} = W^{t+1} - \alpha dW^{t+1}, b^{t+1} = b^{t+1} - \alpha db^{t+1}$   
 }  
 } "1 epoch" pass through training set.

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

$X, Y$

Andrew Ng

미니 배치 기울기 강하를 트레이닝 세트에서 실행하기 위해서는,  $t=1$  에서 5000까지 실행합니다. 각각 1000개씩 가지고 있는 미니 배치가 5000개 있었죠. for loop 내부에서는 무엇을 할 것이냐면,  $x_t, y_t$ 를 이용해서 기울기 강하의 한 단계를 도입할 것입니다. 이것은 마치 트레이닝 세트의 크기가 1000개의 예시가 있는 것과 비슷합니다.

여러분이 별개의 for loop을 1000개의 예시에 갖게 하는 것이 아닌, **벡터화를 통해 1000개의 예시를 한번에 처리하는 것**과 같습니다.

이전에는 여기가 그냥  $X$ 였죠? 하지만 이제 여러분은 전체 트레이닝 세트를 처리하는 것이 아니라, 첫번째 미니 배치만 처리하는 것입니다. 그렇기 때문에 여기가  $x_T$ 가 되죠. 미니 배치  $T$ 를 처리하는 경우에 말이죠.

그러면  $G1$ 의  $Z1$ 의 값을 갖게 되겠죠. 여기는 대문자  $Z$ 죠, 벡터화 도입이기 때문입니다. 이어서  $AL$ 값이 될 때까지 진행하는데요, 여기는  $GL$ 의  $ZL$ 입니다. 그리고 이 값이 예측 값이 되는 것입니다.

다음으로는, **J 비용함수를 산출**합니다. 여기서 이 함수를  $1/1000$ 으로 적을 것인데요, 그 이유는 1000이 여러분의 작은 트레이닝 세트의 크기이기 때문입니다.

여기 제가 쓴 코드는 1 epoch 트레이닝을 한다고 표현하기도 하는데요, epoch이라는 용어는 트레이닝 세트로 1번 통과한다는 뜻입니다. 반면에 **배치 기울기 강하에서는 한번 트레이닝을 통과하는 것이 오로지 한번의 기울기 강하 절차를 밟게 해줍니다.**

# Mini-batch gradient descent

repeat {  
 for  $t = 1, \dots, 5000$  {  
 Forward prop on  $X^{t+1}$ .  
 $Z^{(i)} = W^{(i)} X^{t+1} + b^{(i)}$   
 $A^{(i)} = g^{(i)}(Z^{(i)})$   
 $\vdots$   
 $A^{(n)} = g^{(n)}(Z^{(n)})$  } Vectorized implementation (1000 examples)  
 Compute cost  $J^{t+1} = \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(y^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_i \|W^{(i)}\|_F^2$ .  
 Backprop to compute gradients w.r.t  $J^{t+1}$  (using  $(X^{t+1}, Y^{t+1})$ )  
 $W^{(i)} := W^{(i)} - \alpha dW^{(i)}, b^{(i)} := b^{(i)} - \alpha db^{(i)}$   
 }  
 } "1 epoch" pass through training set.

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

$X, Y$

Andrew Ng

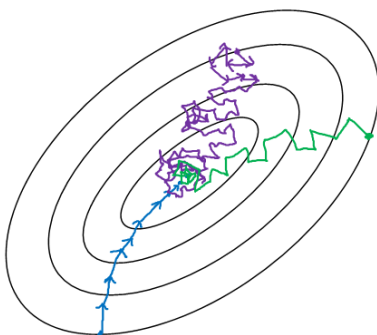
조금 더 noisy한 이유는  $X\{1\}, Y\{1\}$ 이 조금 더 쉬운 미니 배치여서 비용이 조금 더 낮기 때문에 그럴 수 있는데요, 하지만 우연으로  $X\{2\}, Y\{2\}$ 가 어려운 미니 배치일 수도 있죠. mislabel된 예시도 있을 수도 있겠죠, 이런 경우 비용이 조금 더 높겠습니다.

## Choosing your mini-batch size

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

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

In practice: Somewhere in-between 1 and  $m$



Stochastic gradient descent  
 ↳ Lose speedup from vectorization

In-between (mini-batch size not too big/small)

Fastest learning.

- Vectorization. ( $n=1000$ )
- Make passes without processing entire training set.

Batch gradient descent (mini-batch size =  $m$ )

↳ Too long per iteration

Andrew Ng

1. 미니 배치 사이즈를  $m$ 으로 설정하는 것은 단순히 배치 기울기 강하를 줍니다.
2. stochastic 기울기 강하라는 알고리즘에서 각각의 예시는 그들의 미니 배치입니다.

이럴 경우에는, 첫번째 미니 배치를 봅시다. 즉,  $X\{1\}$ ,  $Y\{1\}$ 인데요, 그렇지만 미니 배치 사이즈가 1인 경우, 첫번째 트레이닝 예시 밖에 없는데요, 즉 첫번째 트레이닝 샘플로 기울기 강하를 가져야 합니다. 다음으로 두번째 미니 배치를 봅시다. 두번째는 그냥 두번째 트레이닝 샘플인데요, 이 값에 기울기 강하 step를 적용합니다. 그리고 이어서 세번째 트레이닝 예시도 그렇게 합니다. 하나의 트레이닝 샘플씩 보는 것입니다.

배치 기울기 강하 는 여기쯤에서 시작할텐데요, 이 경우, 낮은 noise와 큰 step를 가질 수 있을 것입니다. 그러면 계속 최소값을 향할 수 있을 것입니다. stochastic 기울기 강하 와는 반대로, 다른 점을 한번 지정해보겠습니다. 그러면 한번의 iteration마다 한개의 트레이닝 예시로 gradient descent를 하는 것인데요,

그러면 만약 미니 배치 사이즈가  $m$ 이 아니고, 1이 아니며, 그 사이 값을 가지려면, 어떤 것을 고를 수 있을까요?

## Choosing your mini-batch size

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

Typical mini-batch sizes:

→  $64, 128, 256, 512$   $\frac{1024}{2^{10}}$   
 $\underbrace{\quad\quad\quad}_{2^6 \quad 2^7 \quad 2^8 \quad 2^9}$

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

Andrew Ng

조금 더 큰 트레이닝 세트의 경우, 전형적인 미니 배치 사이즈는 64에서 512가 대표적인데요, 컴퓨터 메모리가 배치되어 있는 방식과 접속되는 방법에 따라, 가끔씩 그 코드가 2의 지수 값을 가질때 더 빨리 실행됩니다. 64는 그러면 2의 6승이구요, 2의 7승, 2의 8승, 2의 9승, 저는 그래서 자주 미니 배치 사이즈를 어떤값의 2승으로 도입합니다. 이번 비디오에서는 제가 미니 배치 사이즈 1000을 사용했는데요, 만약 정말로 그렇게 하고 싶은 경우엔, 저는 여러분께서 1024를 쓰시는 것을 추천 드립니다. 이 값은 2의 10승이죠. 조금 흔하진 않지만, 미니 배치 사이즈가 1024인게 있긴 합니다. 여기 이 범위의 미니 배치 사이즈가 조금 더 흔하긴 합니다. 마지막 팁으로는,

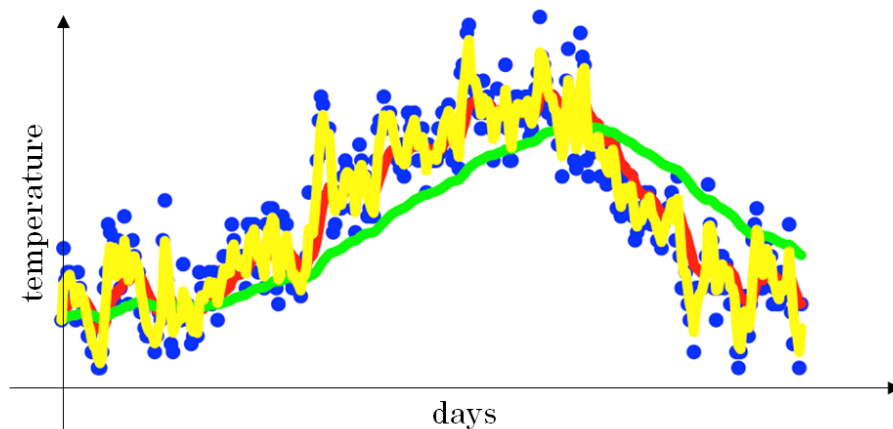
모든  $X\{t\}$ ,  $Y\{t\}$  값이 CPU/GPU 메모리에 들어가게 하도록 하는 것입니다.

## 지수가중 평균법

### Exponentially weighted averages

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

$$\beta = 0.9 \quad 0.98 \quad 0.5$$



Andrew Ng

# Exponentially weighted averages

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

$$\begin{aligned}
 v_{100} &= 0.9v_{99} + 0.1\theta_{100} \\
 v_{99} &= 0.9v_{98} + 0.1\theta_{99} \\
 v_{98} &= 0.9v_{97} + 0.1\theta_{98} \\
 &\dots \\
 \rightarrow v_{100} &= 0.1\theta_{100} + 0.9(0.1\theta_{99} + 0.9(0.1\theta_{98} + 0.9(0.1\theta_{97} + 0.9(0.1\theta_{96} + \dots))) \\
 &= 0.1\theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1(0.9)^2 \theta_{98} + 0.1(0.9)^3 \theta_{97} + 0.1(0.9)^4 \theta_{96} + \dots \\
 &\quad 0.9^{10} \approx 0.35 \approx \frac{1}{2} \quad (1-\epsilon)^{\frac{1}{\epsilon}} = \frac{1}{e} \quad \epsilon = 0.02 \rightarrow 0.98^{50} \approx \frac{1}{e}
 \end{aligned}$$

Andrew Ng

일별 평균 기온을 어떻게 산출하는지 한번 이해해보도록 하겠습니다. 여기는 아까봤던 그 공식인데요, 베타의 값을 0.9로 설정하고, 이 값에 상응하는 몇개의 식을 더 적어보겠습니다. 만약에 T의 값이 0에서 1로 그리고 2에서 3으로 T가 늘어날때의 값을 분석하는 반면에 여기서는 **T의 값이 감소할때의 경우를 적었습니다**. 여기 첫번째 식을 보겠습니다. V100이 실제로는 어떤 의미인지 보겠습니다. V100은 어떻게 되냐면, 여기 2개의 항을 서로 맞바꿔 볼텐데요, 여기는 0.1 곱하기 세타 100에 더하기 0.9 곱하기 전날의 기온을 해줍니다. 그러면 V99는 또 어떻게 되죠? 여기 식에 이 식을 대입해보겠습니다. 이 값은 그냥 0.1 곱하기 세타 99, 그리고 다시 2개의 항을 서로 맞바꿨습니다. 더하기 0.9 곱하기 V98입니다. 그러면 또 V98은 또 어떻게 되죠? 자, 여기서 구할 수 있는데요, 이곳에다가 대입시키면 됩니다.

이런 방법으로 합을 구하면 되고, 세타100에 대한 가중평균치 인데요, 그 해의 100일간의 온도를 계산한 V100 까지 반영하여 현재 시점의 기온을 뜻합니다. 이것은 세타 100과 세타 99 세타 98, 세타 97, 세타 96 등등의 합을 나타내는 값입니다

나중에 제가 설명을 드리겠지만, 여기에 있는 모든 계수들을 더하면 1이 되거나 1과 매우 근접한 값을 갖게됩니다. **편이보정(bias correction)**이라는 값과 이것에 대한 내용은 다음 비디오에서 다루겠지만, 이런 점 때문에 지수 가중 평균이라고 합니다.

베타가 0.9인 경우, 이것은 마치 여러분이 직전 10일 간의 기온만 집중하여 지수 가중 평균을 구하는 것과 같은데요, 약 10일 정도 이후 가중치가 현재 날짜 가중치 3분의 1보다 약간 더 감소하기 때문에 그렇습니다. 반면에 베타의 값이 0.98이었다고 하면, 그러면, 0.98의 몇승을 해야 이 값이 매우 작은 값이 될까요? 알고보니 0.98의 50승이면 대략  $1/e$ 와 비슷한 값이 됩니다. 그렇기 때문에 가중치는 첫 50일에는  $1/e$  보다 큰 값이 되었다가, 그 후엔, 꽤 빠른 속도로 감소할 것입니다.

## Implementing exponentially weighted averages

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$

$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$v_3 = \beta v_2 + (1 - \beta) \theta_3$$

...

```

V_0 := 0
V_0 := beta v + (1-beta) theta_1
V_0 := beta v + (1-beta) theta_2
:

```

---

```

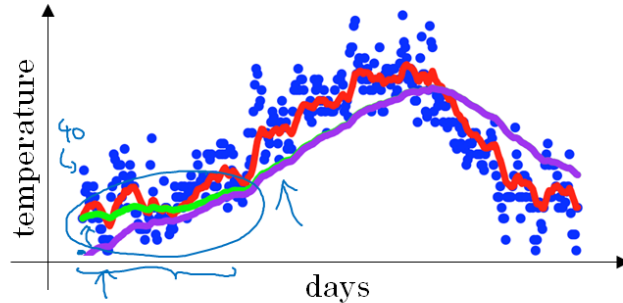
-> V_0 = 0
Repeat {
  Get next theta_t
  V_0 := beta V_0 + (1-beta) theta_t
}

```

Andrew Ng

V를 초기화한 값을 0이라 하고, 첫번째 날은, V의 값을 베타 곱하기 V 더하기 1 빼기 베타 곱하기 세타 1로 설정할 것입니다. 그 다음날에는, V를 업데이트하여 베타 곱하기 V 더하기 1 빼기 베타 곱하기 세타2를 할 것입니다. 이렇게 계속 이어지겠죠. 어떤 것은 이런 V 아래첨자 세타로 표기할텐데요, 세타인 파라미터에 대해 지수 가중 평균을 구하고 있다는 의미를 나타냅니다.

# Bias correction



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

$$v_0 = 0$$

$$v_1 = 0.98 v_0 + 0.02 \theta_1$$

$$v_2 = 0.98 v_1 + 0.02 \theta_2$$

$$= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2$$

$$= 0.0196 \theta_1 + 0.02 \theta_2$$

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

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

$$\frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

Andrew Ng

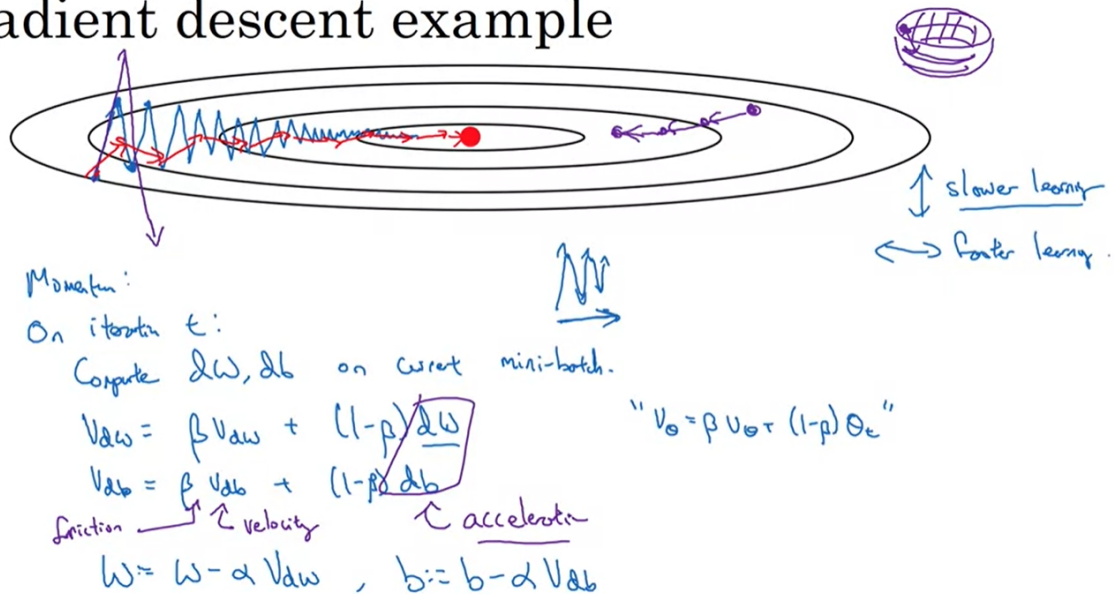
이 그래프는 0.98인 경우입니다. 베타가 0.98일때 말이죠. 대신 이렇게 생긴 보라색 곡선을 갖게됩니다. 여러분도 보시겠지만 보라색 곡선은 굉장히 낮은데서 시작합니다. 이것이 끼치는 영향을 보겠습니다. **우리가 이동평균을 구현할 때,  $v_0 = 0$ 으로 초기화합니다.** 그리고  $v_1 = 0.98 v_0 + 0.02$  제타 1 이 됩니다. 그러나  $v_0$ 은 0이기 때문에 이 항은 그냥 사라지게 됩니다.  $v_1$ 은 그러면 0.02 곱하기 제타1입니다. 그렇기 때문에 첫째 날의 온도가, 예를 들어, 화씨 40도 인 경우,  $v_1$ 은 0.02 곱하기 40 이므로 0.8 이 됩니다. 그래서, 훨씬 더 작은 값이 여기서 나오게 됩니다. 그러면 첫째 날 온도로 좋은 추정치가 아니죠.  $v_2$ 는 0.98 곱하기  $v_1$ , 더하기 0.02 곱하기 제타 2 이고,  $v_1$ 은 이것을 아래로 대입해서 곱하면  $v_2$ 는 0.98 곱하기 0.02 곱하기 제타 1 더하기 0.02 곱하기 제타2가 됩니다. 그러면 0.0196 제타 1 더하기 0.02 제타 2 이고요. 그리고 다시 제타1과 제타2가 양수라 가정하면,  $v_2$ 는 제타1이나 제타2보다 훨씬 작은 값이 될 것 입니다. **그래서,  $v_2$ 는 그 해의 첫 이틀간 기온에 대한 좋은 추정치가 아닙니다.**

이 추정치를 더 낮게 보완하는 방법이 있는데요, 더 정확하게 만들어 주는 것인데. 특히 추정 초반 부분에 대한 것입니다.  **$v_t$  대신에  $v_t$  나누기 1 빼기 베타의  $t$  승을 사용하는 것인데,**  $t$ 는 이 시점의 현재 데이터 입니다. 구체적인 예를 살펴보죠.  $t$ 의 값이 2일 때, 1 빼기 베타의  $t$ 승은  $1 - 0.98$  제곱이고, 이 값은 0.0396입니다. 그러므로 두째 날의 기온 추정치는  $v_2$  나누기 0.0396이며, 이 값은 0.0196 곱하기 제타1 더하기 0.02 곱하기 제타 2가 됩니다. 보시면, 여기 이 둘에 분모로 0.0396으로 하고, 이것이 제타1과 제타2의 가중평균이 되서, 이 바이어스를 제거됩니다. 그러므로  $t$ 의 값이 커지면서 베타의  $t$ 승은 0으로 수렴하고,  $t$  값이 충분히 커



지면, 바이어스 보정 영향은 없게 됩니다. 그러므로  $t$ 의 값이 큰 경우, 보라색 선과 초록색 선은 훨씬 많이 겹치게 됩니다. 이 추정치의 초반부에는 추정치가 서서히 맞춰 들어 갈 때, 바이어스 보정을 통해서 기존의 더 나은 추정치를 얻을 수 있고요. 바이어스 보정으로 보라색 선이 초록색 선으로 가게 됩니다.

## Gradient descent example



Andrew Ng

# Implementation details

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

On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

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

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

$$\frac{v_{dw}}{1 - \beta^t}$$

Hyperparameters:  $\alpha, \beta$

$$\beta = 0.9$$

average over last  $\approx 10$  gradients

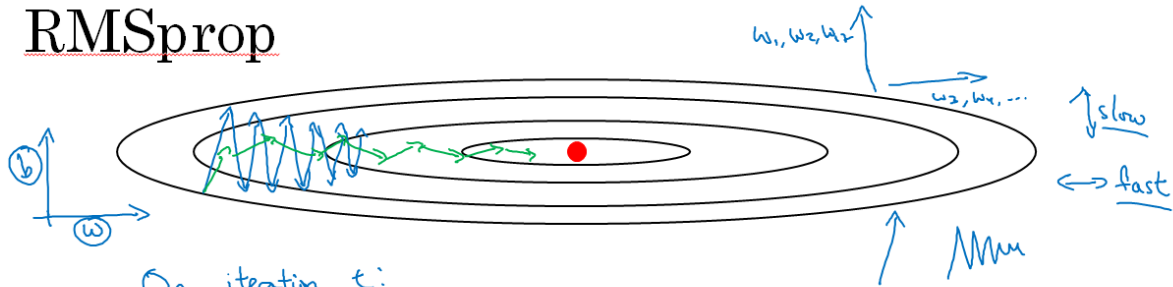
Andrew Ng

마지막으로 여러분이 기울기 강하 와 모멘텀에 대한 학술을 읽으시면, 이 항이 생략되는 것을 보실텐데요 1 빼기 베타 부분말이죠. 그렇게해서  $vdW = \text{베타 } vdw + dW$ 로 남게됩니다. 여기 보라색 버전을 쓰면서 나타나는 최종 효과는  $vdW$ 가 1 빼기 1베타로 스케일 된다는 것입니다. 또는, 더 정확히 얘기하면 1 나누기 1-베타로 말입니다. 그러므로 여러분이 기울기 강하 갱신을 할때 알파의 값은 그에 상응하는 1 나누기 1 빼기 베타로 변해야 합니다.

## RMSprop

기울기 강하를 도입하는 경우엔, 세로로 아주 큰 변동이 일어날 수 있습니다. 가로줄로 나아가려고 할때 말이죠.

# RMSprop



On iteration  $t$ :

Compute  $dW, db$  on current mini-batch

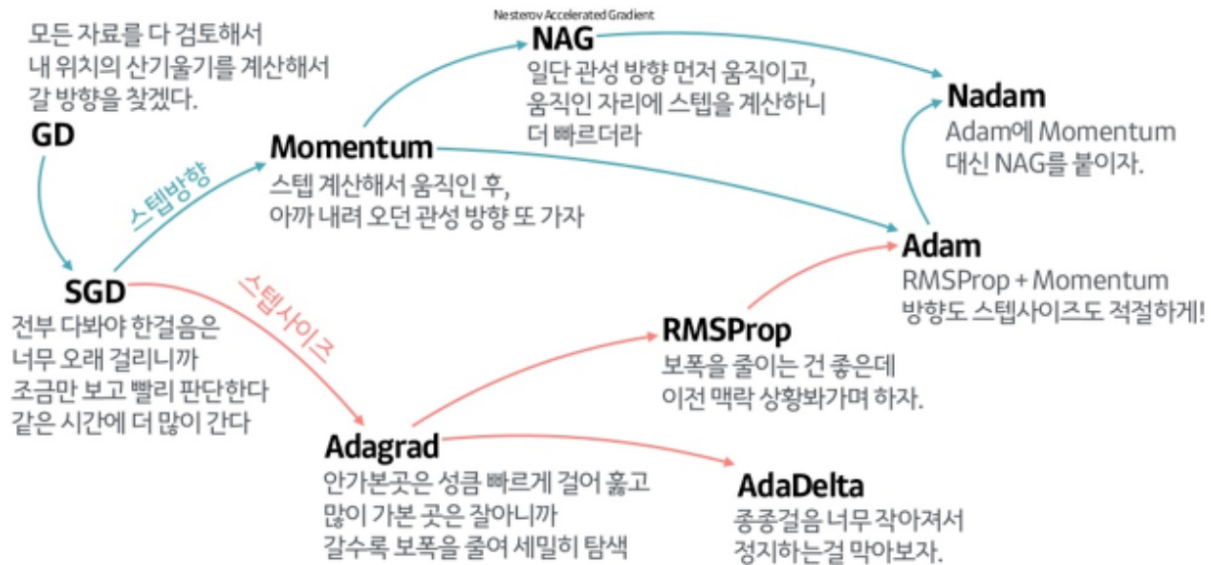
$$S_{aw} = \beta_2 S_{aw} + (1 - \beta_2) \underbrace{dW^2}_{\text{element-wise}} \leftarrow \text{small}$$

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

$$w := w - \alpha \frac{dW}{\sqrt{S_{aw} + \epsilon}} \quad b := b - \alpha \frac{db}{\sqrt{S_{ab} + \epsilon}}$$

$$\epsilon = 10^{-8}$$

Andrew Ng



RMSprop

- AdaGrad는 스텝이 많이 진행되면 누적치  $h_n$ 이 너무 커져서 학습률이 너무 작아져 학습이 거의 되지 않는 문제가 있음
- 이를 보완하기 위해 RMSProp은 Adagrad와 달리 기울기를 단순 누적 x
- AdaGrad보다 최근 값을 더 잘 반영하기 위해 최근 값과 이전 값에 각각 가중치를 주어 계산하는 방법
- **지수 이동 평균** 사용
- 쉽게 말하면, 내분인데, **더 중요한 변수**에 내분점을 **가까이** 찍어서 가중치를 더 주는 것
- $\gamma$ 가 **크다** → 빨간점에 더 가까이 → **과거** 중요하게 생각
- $\gamma$ 가 **작다** →  $(1 - \gamma)$ 가 크다 → 파란점에 더 가까이 → **현재** 중요하게 생각

$$h_n = \gamma h_{n-1} + (1 - \gamma) \nabla f(\mathbf{x}_n) \odot \nabla f(\mathbf{x}_n), \quad h_{-1} = 0$$



Adam : momentum + RMSprop

# 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

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

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

$$\text{vhat} = \text{np.array}([.9, 0.2, 0.1, .4, .9])$$

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

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

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

Andrew Ng

미분을 계산할텐데요.  $dw$ 와  $db$ 를 현재 미니 배치를 이용해 계산하고 보통, 미니 배치 기울기 강하를 이용하면 되고요.

그리고 모멘텀, 지수 가중 평균을 쓸텐데요.  $V_{dw} = \beta$ 인데, 이제부터는 이 값을  $\beta_1$ 이라고 쓰고, **RmsProp의 하이퍼 파라미터  $\beta$ 는  $\beta_2$ 로 구분해서** 사용하겠습니다.

자, 이것은 모멘텀을 구현할 때 했던 것과 완전히 동일한데요 유일한 차이는 하이퍼 파라미터  $\beta$  대신에  $\beta_1$ 이라고 부른 점입니다. 유사하게  $V_{db}$ 는 이렇게, 1 빼기  $\beta_1$  곱하기  $db$ 죠. 그리고 RmsProp도 업데이트 하는데요. 이제 하이퍼 파라미터는  $\beta_2$ 죠. 여기는 플러스 1 빼기  $\beta_2$   $dw^2$  여기서 제공은 미분  $dw$ 를 원소별로 제공하는 것이고요.

일반적으로 Adam 구현에는, **편향보정도** 같이 합니다. 그러므로 **v corrected** 를 사용할텐데, Corrected(보정된)는 편향 보정된 것을 뜻 합니다.  $dw$ 는  $v_{dw}$  나누기 1 빼기  $\beta_1$ 의  $t$ 승,  $t$ 회째 반복 루프를 진행한 경우에 말이죠. 비슷하게,  $v_{db}$  corrected는  $v_{db}$  나누기 1 빼기  $\beta_1$ 의  $t$ 승입니다. 그리고 또 비슷하게,  $S$ 에도 편향보정을 적용해서 즉,  $sdw$  나누기 1 빼기  $\beta_2$ 의  $T$ 승, 그리고  $sdb$  corrected는  $sdb$  나누기 1 빼기  $\beta_2$ 의  $T$ 승입니다. 마지막으로 업데이트를 진행합니다. 그러면  $W$ 는  $W$  빼기 알파 곱하기, 만약 모멘텀만 사용하는 경우  $V_{dw}$ 를 사용하거나,  $V_{dw}$  corrected가 될 수도 있겠죠. 이제 RmsProp에 해당하는 부분도 추가하는데요. 그래서,  $S_{dw}$  corrected + 애플론의 루트로 나눠줍니다. 마찬가지로, 비슷한 공식으로  $b$  또한 업데이트 될텐데요.  $V_{db}$  corrected 나누기 루트  $S$  corrected  $db$  더하기 애플론입니다. 이 알고리즘은 기울기 강하에 모멘텀 효과와 함께 RmsProp을 같이 결합한 것입니다.

## Hyperparameters choice:

→  $\alpha$  : needs to be tune  
→  $\beta_1$  : 0.9 → (dw)  
→  $\beta_2$  : 0.999 → (dw<sup>2</sup>)  
→  $\epsilon$  :  $10^{-8}$

Adam: Adaptive moment estimation



Adam Coates

Andrew Ng

이 알고리즘은 몇개의 하이퍼 파라미터들이 있는데요. 러닝 속도인 하이퍼 파라미터 알파 $\alpha$ 는 여전히 중요하고, 보통 튜닝이 필요합니다. 그러므로 여러분이 다양한 범위의 값을 시도해서 어떤 값이 적합한지 찾아야 합니다.  $\beta_1$ 의 가장 흔한 설정 값은 0.9입니다. 이것은 이동 평균값인데요.  $dw$ 의 가중평균값인데요, 이것은 모멘텀 같은 성분이고, 하이퍼 파라미터  $\beta_2$ 는, Adam 논문의 저자, Adam 알고리즘 발명자는 0.999를 권장하였습니다. 이것은  $dw^2$ 과 또,  $db^2$ 의 이동평균 (지수 가중 평균)들을 계산하는 것입니다. 그리고, 앵실론 $\epsilon$  값 선택은 그리 중요하지 않습니다.

Adam은 "적응 모멘트 추정" (Adaptive Moment Estimation)을 뜻합니다.  **$\beta_1$ 은 미분의 평균값을 산출하는데요. 이것을 1차 모멘트라고 합니다.** 그리고  **$\beta_2$ 는  $dW^2$ 의 지수이동평균에 사용되는데, 이것은 2차 모멘트라고 합니다.** 이런 이유에서 adaptive moment estimation이라는 이름이 나오게 되었습니다.

# Learning rate decay

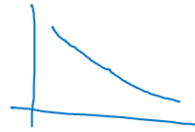
1 epoch = 1 pass through data.

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

Epoch	$\alpha$
1	0.1
2	0.67
3	0.5
4	0.4
...	...



$\alpha_0 = 0.2$   
decay-rate = 1



Andrew Ng

학습 알고리즘을 가속화시키는데 도움이 될 수 있는 것 중 하나는 러닝 레이트를 시간이 가면서 천천히 줄여나가는 것입니다. 이것을 learning rate decay (학습 속도 감쇠법)라고 합니다.

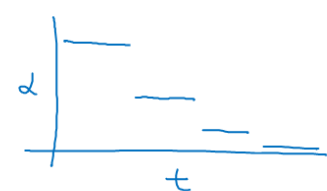
여러 에폭(epoch)을 진행해 보면 즉, 데이터에 여러번 패스를 돌리면  $\alpha_0$  는 0.2 이고, 감쇠비(decay-rate)는 1이면, 그럼 첫번 에폭(epoch)에서  $\alpha$ 는  $1 / (1 + 1 * \alpha_0)$ 입니다. 학습 속도는 0.1이 되겠죠. 이 공식에 넣어보면 나오고요. 감쇠비(decay-rate)는 1이고, 에폭 값(epoch-num)은 1일 때 말이죠. 두번째 에폭(epoch)에서는 러닝 속도가 0.67로 감쇠되었습니다.

세번째에서는 0.5, 네번째에서는 0.4. 등등.

학습 속도는 점진적으로 작아지는 것이 감이 오죠.

## Other learning rate decay methods

Formula

$$\left\{ \begin{array}{l} \alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \quad - \text{exponentially decay.} \\ \alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \cdot \alpha_0 \\ \text{discrete staircase} \end{array} \right.$$


Manual decay.

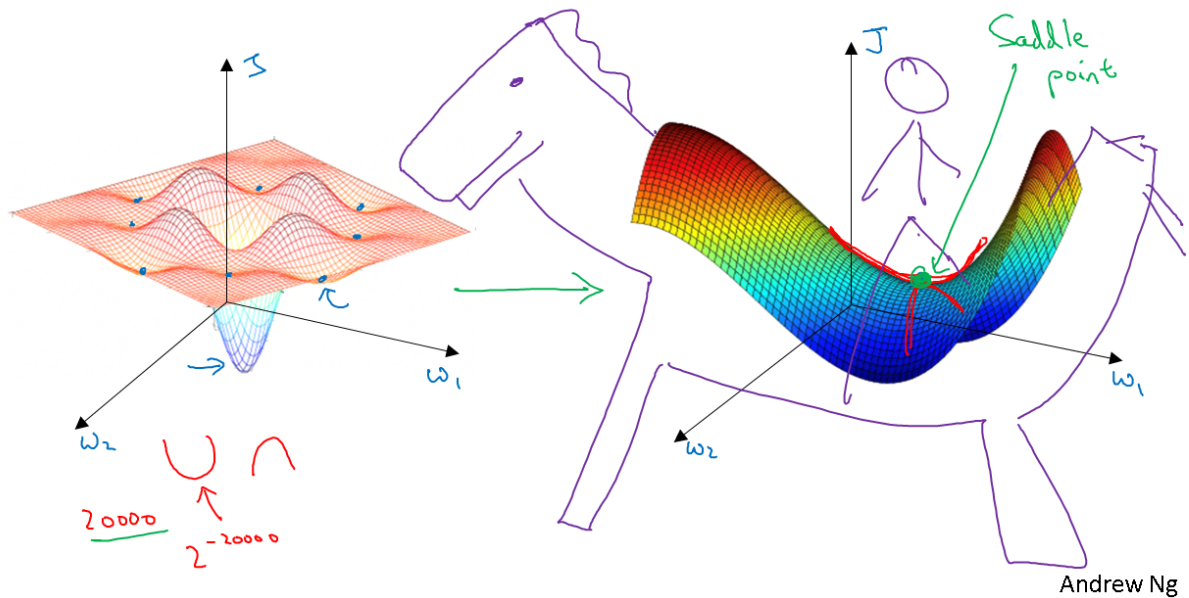
Andrew Ng

사람들이 쓰는 또 다른 공식은  $\alpha = \text{어떤 상수 } k / \text{루트 에폭값에 곱하기 } \alpha_0$ . 또는, 또 하나의 하이퍼 파라미터죠. 어떤 상수  $k$ 에 나누기 루트 미니 배치의 수  $t$  곱하기  $\alpha_0$  입니다.

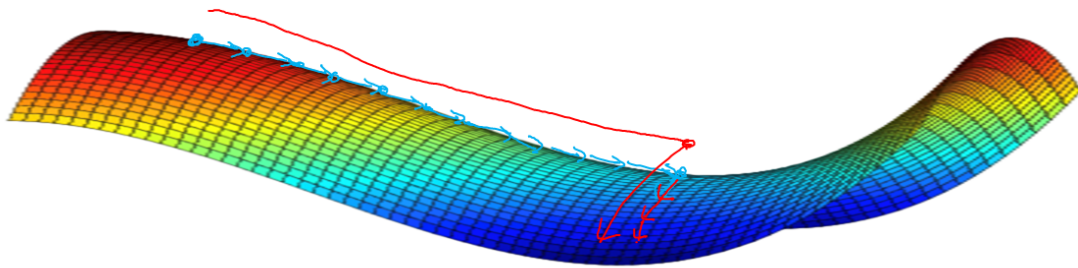
가끔 사람들이 하는 한 가지 방법은 **수작업식 감쇠(manual decay)**입니다.



# Local optima in neural networks



## Problem of plateaus



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

Andrew Ng

국소 최적값이 문제가 되지 않는다면, 어떤 것이 문제가 될까요? plateau 가 러닝속도를 저하시킬 수 있는데요, plateau는 함수 기울기의 값이 0에 근접한 긴 범위를 이야기 합니다.

momentum 또는 RmsProp 또는 Adam 과 같은 알고리즘이 도움을 줄 수 있는 부분입니다. 이런 plateau와 같은 시나리오 경우, 가장 정교한 observation 알고리즘인 Adam과 같은 알고리즘이 plateau를 빠져나오는 속도를 높여줄 수 있습니다.