

# Deep Learning

## 11 Training DNNs II

Dr. Konda Reddy Mopuri  
Dept. of AI, IIT Hyderabad  
Jan-May 2023

# Issues with SGD

- DNNs are trained via SGD:  $w_{t+1} = w_t - \eta \cdot \nabla_w J(w)$

# Issues with SGD

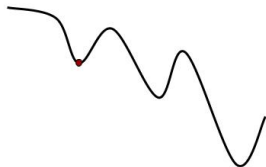
- DNNs are trained via SGD:  $w_{t+1} = w_t - \eta \cdot \nabla_w J(w)$
- Loss is a high dimensional function

# Issues with SGD

- DNNs are trained via SGD:  $w_{t+1} = w_t - \eta \cdot \nabla_w J(w)$
- Loss is a high dimensional function
  - May have local minima

# Issues with SGD

- DNNs are trained via SGD:  $w_{t+1} = w_t - \eta \cdot \nabla_w J(w)$
- Loss is a high dimensional function
  - May have local minima
  - May have saddle points



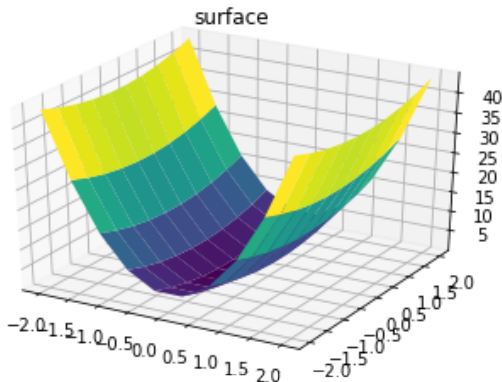
Stuck at a local minimum



Stuck at a saddle point

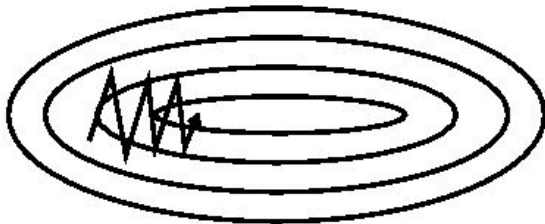
# Issues with SGD

- DNNs are trained via SGD:  $w_{t+1} = w_t - \eta \cdot \nabla_w J(w)$
- Loss is a high dimensional function
  - May vary swiftly in direction and slowly in the other



# Issues with SGD

- SGD leads to jitter along the deep dimension and slow progress along the shallow one



---

Figure credits: Sebastian Ruder

## SGD+Momentum

### SGD

$$w_{t+1} = w_t - \eta \cdot \nabla_w J(w)$$

$$v_0 = 0$$

$$v_{t+1} = \rho \cdot v_t + \nabla_w J(w)$$

$$w_{t+1} = w_t - \eta \cdot v_{t+1}$$

---

I Sutskever et al., ICML 2013



## SGD+Momentum

### SGD

$$w_{t+1} = w_t - \eta \cdot \nabla_w J(w)$$

$$v_0 = 0$$

$$v_{t+1} = \rho \cdot v_t + \nabla_w J(w)$$

$$w_{t+1} = w_t - \eta \cdot v_{t+1}$$

- Aggregates velocity: exponential moving average over gradients

## SGD+Momentum

### SGD

$$w_{t+1} = w_t - \eta \cdot \nabla_w J(w)$$

$$v_0 = 0$$

$$v_{t+1} = \rho \cdot v_t + \nabla_w J(w)$$

$$w_{t+1} = w_t - \eta \cdot v_{t+1}$$

- Aggregates velocity: exponential moving average over gradients
- $\rho$  is the friction (typically set to 0.9 or 0.99)

---

I Sutskever et al., ICML 2013

# SGD+Momentum

## SGD+Momentum

### SGD

$$w_{t+1} = w_t - \eta \cdot \nabla_w J(w)$$

```
for i in range(num_iters):  
    → dw = grad(J, W, x, y)  
    → w- = η · dw
```

$$v_0 = 0$$

$$v_{t+1} = \rho \cdot v_t + \nabla_w J(w)$$

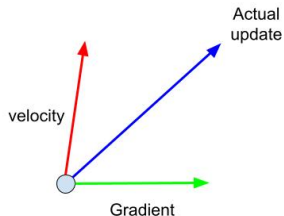
$$w_{t+1} = w_t - \eta \cdot v_{t+1}$$

```
v_0 = 0  
for i in range(num_iters):  
    → dw = grad(J, W, x, y)  
    → v = ρ · v + dw  
    → w- = η · v
```

---

I Sutskever et al., ICML 2013

# SGD+Momentum



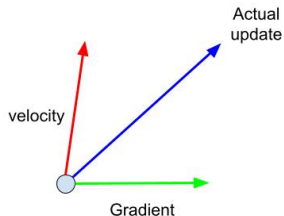
Momentum Update

① How can momentum help?

---

I Sutskever et al., ICML 2013

# SGD+Momentum



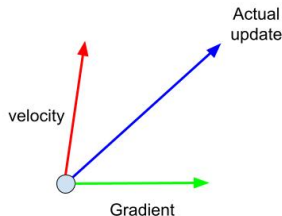
Momentum Update

- ① How can momentum help?
  - Optimization proceeds even at the local minimum or saddle point (because of the accumulated velocity)

---

I Sutskever et al., ICML 2013

# SGD+Momentum



Momentum Update

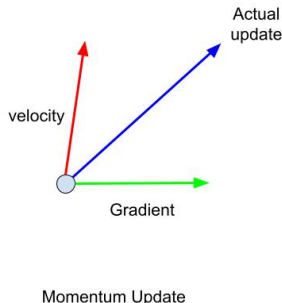
## ① How can momentum help?

- Optimization proceeds even at the local minimum or saddle point (because of the accumulated velocity)
- Jitter is reduced in ravine like loss surfaces

---

I Sutskever et al., ICML 2013

# SGD+Momentum



## ① How can momentum help?

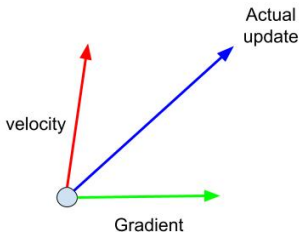
- Optimization proceeds even at the local minimum or saddle point (because of the accumulated velocity)
- Jitter is reduced in ravine like loss surfaces
- Updates are more smoothed out (less noisy because of the exponential averaging)

---

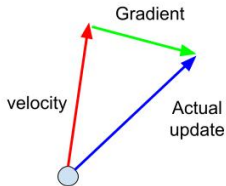
I Sutskever et al., ICML 2013

# Nesterov Momentum

- ① Look ahead with the velocity, then take a step in the gradient's direction



Momentum Update



Nesterov Momentum

---

I Sutskever et al., ICML 2013



# Nesterov Momentum

$$v_0 = 0$$

```
for i in range(num_iters):
```

```
→ dw = grad(J, W +  $\rho \cdot v$ , x, y)
```

```
→  $v = \rho \cdot v + dw$ 
```

```
→  $w = w + \eta \cdot v$ 
```

---

I Sutskever et al., ICML 2013

- ① Adaptive (or, per-parameter) learning rates are introduced

---

Duchi et al. 2011, JMLR

- ① Adaptive (or, per-parameter) learning rates are introduced
- ② Parameter-wise scaling of the learning rate by the aggregated gradient

---

Duchi et al. 2011, JMLR

```
grad_sq = 0
for i in range(max_iters):
    → dw = →grad(J,w,x,y)
    →grad_sq += dw
    →  $w = w - \eta \cdot dw / (\text{sqrt}(\text{grad\_sq}) + \epsilon)$ 
```

- Optimization progress along the steep directions is attenuated

---

Duchi et al. 2011, JMLR

```
grad_sq = 0
for i in range(max_iters):
    → dw = →grad(J,w,x,y)
    →grad_sq += dw
    →  $w = w - \eta \cdot dw / (\text{sqrt}(\text{grad\_sq}) + \epsilon)$ 
```

- Optimization progress along the steep directions is attenuated
- Along the flat directions is accelerated

---

Duchi et al. 2011, JMLR

- ① If Ada Grad is run for too long
  - the gradients accumulate to a big value
  - $\rightarrow$  update becomes too small (or, learning rate is reduced continuously)

- ① If Ada Grad is run for too long
  - the gradients accumulate to a big value
  - $\rightarrow$  update becomes too small (or, learning rate is reduced continuously)
- ② RMS prop (a leaky version of Ada Grad) addresses this using a friction coefficient ( $\rho$ )

```
grad_sq = 0
for i in range(max_iters):
    → dw = →grad(J,w,x,y)
    →grad_sq =  $\rho \cdot \text{grad\_sq} + (1 - \rho) \cdot dw$ 
    →  $w- = \eta \cdot dw / (\text{sqrt}(\text{grad\_sq}) + \epsilon)$ 
```



- ① Inculcates both the good things: momentum and the adaptive learning rates

Adam = RMSProp + Momentum

- ① Inculcates both the good things: momentum and the adaptive learning rates

Adam = RMSProp + Momentum

- ②  $m1 = 0$

$m2 = 0$

```
for i in range(max_iters):
```

```
→ dw = grad(J,w,x,y)
```

```
→  $m1 = \beta_1 \cdot m1 + (1 - \beta_1) \cdot dw$ 
```

```
→  $m2 = \beta_2 \cdot m2 + (1 - \beta_2) \cdot dw^2$ 
```

```
→  $w- = \eta \cdot m1 / (\text{sqrt}(m2) + \epsilon)$ 
```

```
①  $m1 = 0$   
 $m2 = 0$   
for i in range(max_iters):  
    →  $dw = \text{grad}(J, w, x, y)$   
    →  $m1 = \beta_1 \cdot m1 + (1 - \beta_1) \cdot dw$   
    →  $m2 = \beta_2 \cdot m2 + (1 - \beta_2) \cdot dw^2$   
    →  $w- = \eta \cdot m1 / (\text{sqrt}(m2) + \epsilon)$ 
```

- ①  $m1 = 0$   
 $m2 = 0$   
for i in range(max\_iters):  
→  $dw = \text{grad}(J, w, x, y)$   
→  $m1 = \beta_1 \cdot m1 + (1 - \beta_1) \cdot dw$   
→  $m2 = \beta_2 \cdot m2 + (1 - \beta_2) \cdot dw^2$   
→  $w- = \eta \cdot m1 / (\text{sqrt}(m2) + \epsilon)$
- ② Bias correction is performed (since the estimates start from 0)

①  $m1 = 0$

$$m2 = 0$$

for i in range(max\_iters):

→  $dw = \text{grad}(J, w, x, y)$

→  $m1 = \beta_1 \cdot m1 + (1 - \beta_1) \cdot dw$

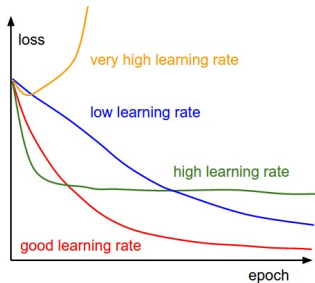
→  $m2 = \beta_2 \cdot m2 + (1 - \beta_2) \cdot dw^2$

→  $w- = \eta \cdot m1 / (\text{sqrt}(m2) + \epsilon)$

② Bias correction is performed (since the estimates start from 0)

③ Adam works well in practice (mostly with a fixed set of values for the hyper-params)

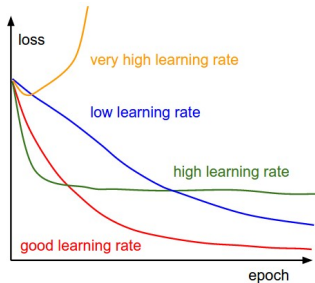
# Learning rate ( $lr$ )



● What  $lr$  to use?

Figure credits: CS231n-Stanford

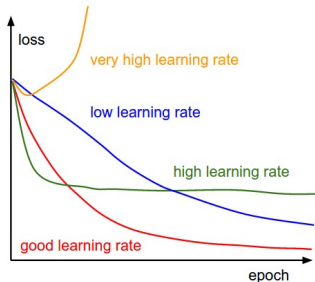
# Learning rate ( $lr$ )



- What  $lr$  to use?
- Different  $lr$  at different stages of the training!

Figure credits: CS231n-Stanford

# Learning rate ( $lr$ )

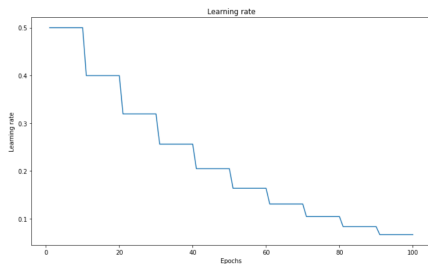


- What  $lr$  to use?
- Different  $lr$  at different stages of the training!
- Start with high  $lr$  and reduce it with time

Figure credits: CS231n-Stanford



# Learning Rate decay: Step

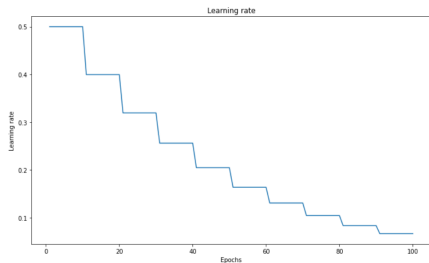


- 1 Reduce the  $lr$  after regular intervals

---

Figure credits: Katherine Li

# Learning Rate decay: Step

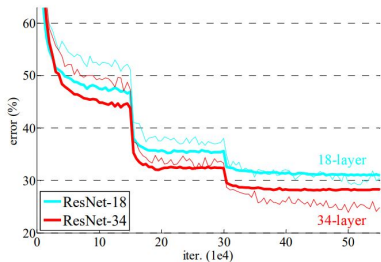


- 1 Reduce the  $lr$  after regular intervals
- 2 E.g. after every 30 epochs,  $\eta^* = 0.1 \cdot \eta$

---

Figure credits: Katherine Li

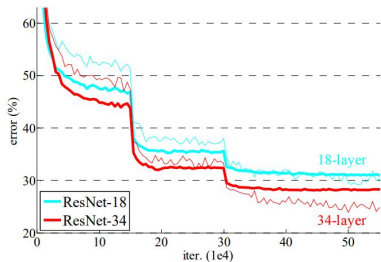
# Learning Rate decay: Step



- ① Characteristic loss curve: different phases for 'stage'

Figure credits: Kaiming He et al. 2015, ResNets

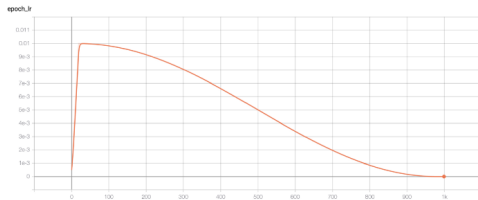
# Learning Rate decay: Step



- 1 Characteristic loss curve: different phases for 'stage'
- 2 Issues: annoying hyper-params (when to reduce, by how much, etc.)

Figure credits: Kaiming He et al. 2015, ResNets

# Learning Rate decay: Cosine



- ① Reduces the  $lr$  continuously
- $$\eta_t = \frac{1}{2}\eta_0(1 + \cos(t\pi/T))$$

Figure credits: Sebastian Correa and Medium.com

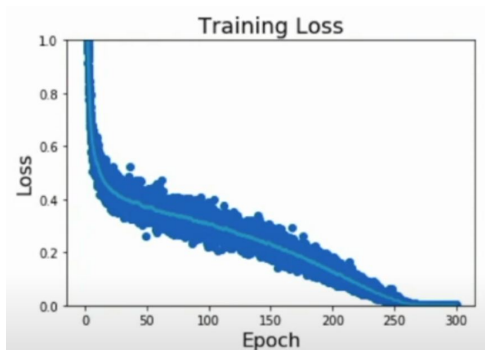
# Learning Rate decay: Cosine



- 1 Reduces the  $lr$  continuously  
$$\eta_t = \frac{1}{2}\eta_0(1 + \cos(t\pi/T))$$
- 2 Less number of hyper-parameters

Figure credits: Sebastian Correa and Medium.com

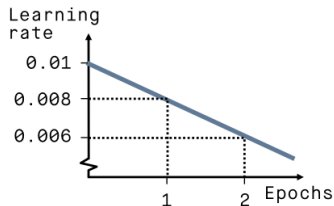
# Learning Rate decay: Cosine



- ① Training longer tends to work, but initial  $lr$  is still a tricky one

Figure credits: Dr Justin Johnson, U Michigan

# Learning Rate decay: Linear

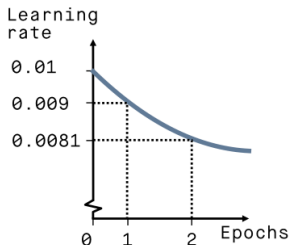


$$\textcircled{1} \quad \eta_t = \eta_0(1 - t/T)$$

Figure credits: [peltarion.com](https://peltarion.com)



# Learning Rate decay: Exponential



$$\textcircled{1} \eta_t = \eta_0 \cdot (1 - \alpha/100)^t$$

Figure credits: [peltarion.com](https://peltarion.com)

# Learning Rate decay: Constant $lr$

- ① No change in the learning rate

$$\eta_t = \eta_0$$

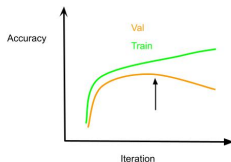
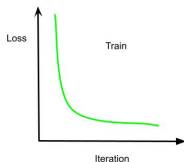
# Learning Rate decay: Constant $lr$

- ① No change in the learning rate

$$\eta_t = \eta_0$$

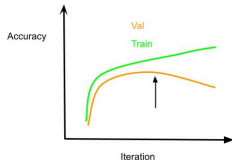
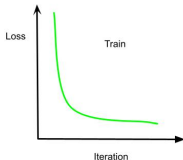
- ② Works for prototyping of ideas (other schedules may be better for squeezing in those 1-2% of gains in the performance)

# Early stopping



- 1 Train as long as the validation performance improves (Stop when it deteriorates)

# Early stopping



- ① Train as long as the validation performance improves (Stop when it deteriorates)
- ② Practice: train for a long number of epochs, saving the intermediate snapshots regularly, pick the one with the best val performance!

# Good training practices

- Observe the initial loss value (if it is as expected or presence of bugs!)

# Good training practices

- Observe the initial loss value (if it is as expected or presence of bugs!)
- One may try to overfit to a very small subset to ensure the basic things are in place

# Good training practices

- Observe the initial loss value (if it is as expected or presence of bugs!)
- One may try to overfit to a very small subset to ensure the basic things are in place
- Monitor the learning curves (tell us if poor initialization or over/under/right-fitting)



# Good training practices

- Observe the initial loss value (if it is as expected or presence of bugs!)
- One may try to overfit to a very small subset to ensure the basic things are in place
- Monitor the learning curves (tell us if poor initialization or over/under/right-fitting)
- Use frameworks' (or fora) help for observing the learning dynamics (e.g. Tensorboard)

# Model Ensembles

- Train multiple models independently and take average inference during testing

# Model Ensembles

- Train multiple models independently and take average inference during testing
- Generally results in slight performance improvements

# Model Ensembles

- The experts can be different snapshots of the same model from training

# Model Ensembles

- The experts can be different snapshots of the same model from training
- E.g. trained with a periodic  $lr$  scheduling

- Moving average of parameters for testing (Polyak Averaging)  
for  $i$  in range(max\_iters):  
→  $dw = \text{grad}(J, w, x, y)$   
→  $w \leftarrow w - \eta \cdot dw$   
→  $w_{\text{test}} = 0.95 \cdot w_{\text{test}} + 0.05 \cdot w$

# Transfer learning: Pretrained features

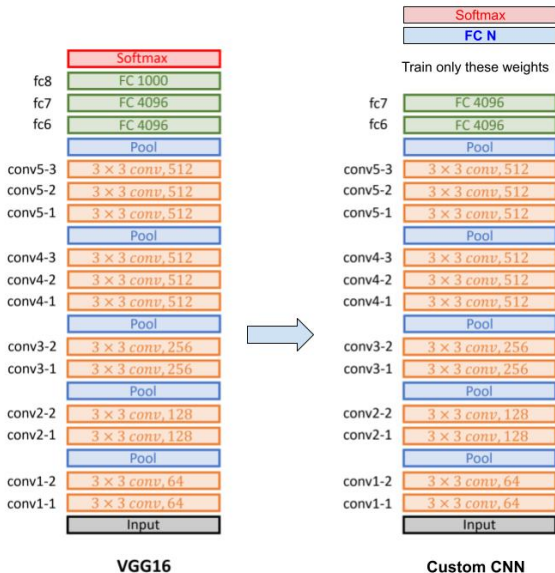
- Sometimes, we may get away with lesser training data!

# Transfer learning: Pretrained features

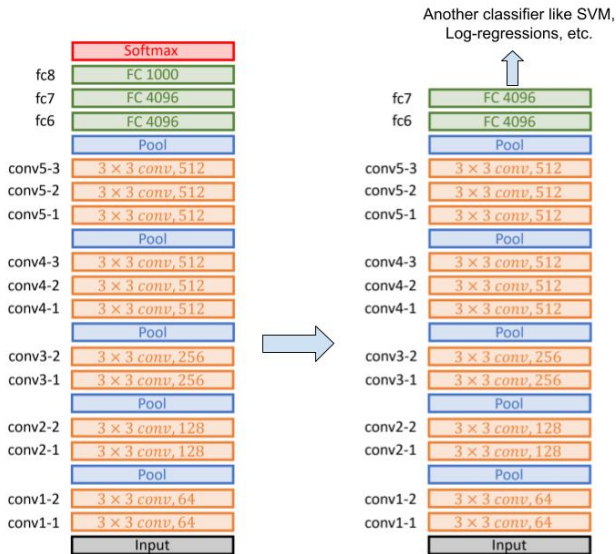
- Sometimes, we may get away with lesser training data!
- Take a DNN trained on a huge training data (task), use it as a feature extractor!!



# Transfer learning: Pretrained features

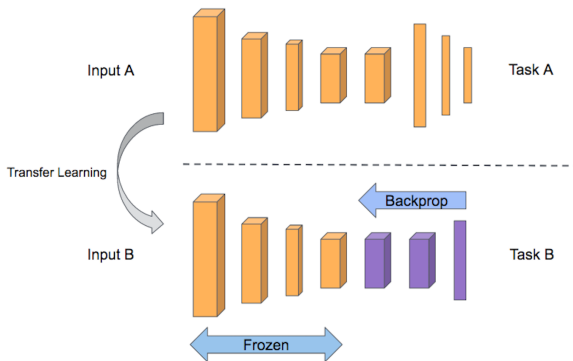


# Transfer learning: Pretrained features



VGG16

# Transfer learning: Pretrained features and Finetuning



Some tips: may have to use smaller learning rate for the transferred layers, start with feature extraction then do finetuning, lower layers might be frozen, etc.

Figure credits: [Giang Tran and Medium.com](#)

# Convex function

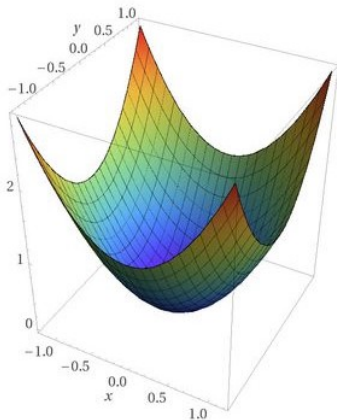


Figure credits: Paperspace blog

# Level sets and ravine

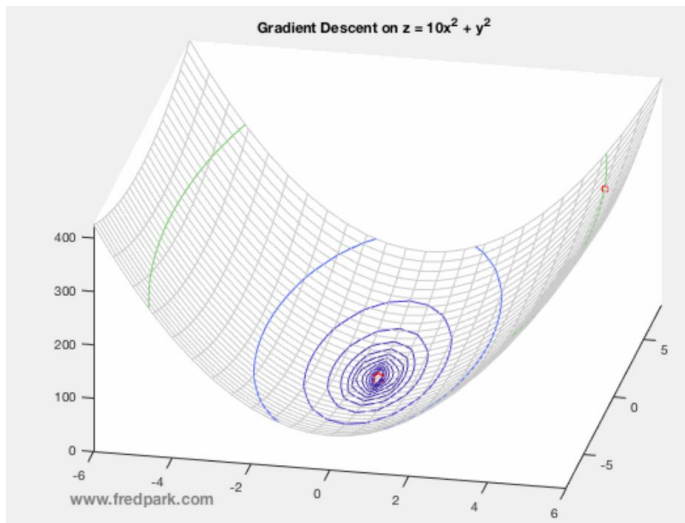


Figure credits: fredpark.com