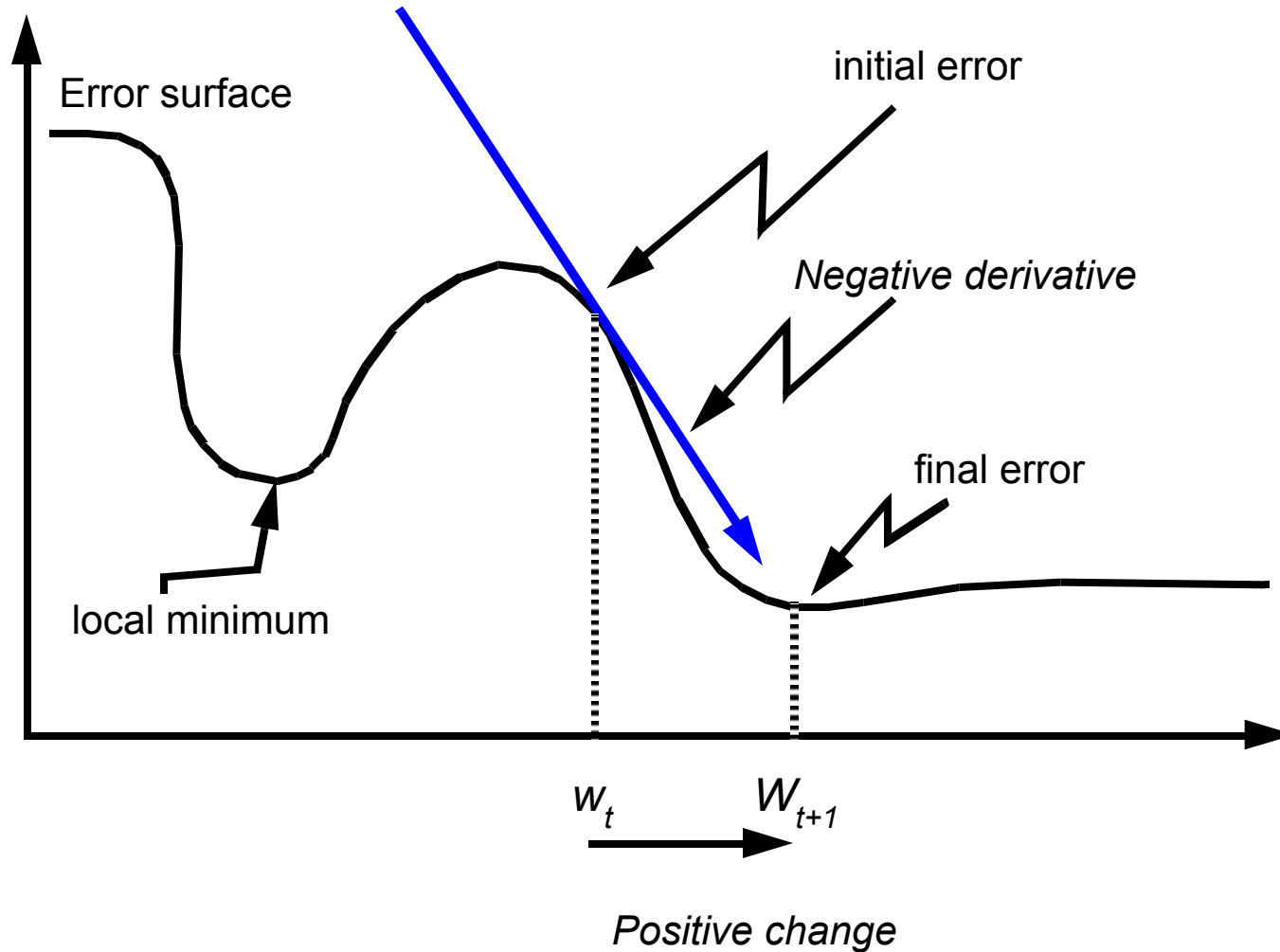


# DEEP LEARNING GRADIENT DESCENT & OPTIMIZERS

Dr. Omri Allouche  
Y-Data Deep Learning Course

# Minimizing the Error



# Gradient Descent

Suppose we have a **scalar function**

We want to find a local minimum.

Assume our current weight is  $w$

$$f(w) : \mathbb{R} \rightarrow \mathbb{R}$$

GRADIENT DESCENT RULE:

$$w_{t+1} \leftarrow w_t - \eta \frac{\partial}{\partial w} f(w_t)$$

$\eta$  is called the LEARNING RATE.

It's a small positive number, e.g.  $\eta = 0.05$

# Gradient Descent in “M” Dimensions

Given  $f(\mathbf{w}) : \mathbb{R}^m \rightarrow \mathbb{R}$

$$\nabla f(\mathbf{w}) = \begin{pmatrix} \frac{\partial}{\partial w_1} f(\mathbf{w}) \\ \vdots \\ \frac{\partial}{\partial w_m} f(\mathbf{w}) \end{pmatrix} \text{ points in direction of steepest ascent.}$$

$|\nabla f(\mathbf{w})|$  is the gradient in that direction

GRADIENT DESCENT RULE:  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \nabla f(\mathbf{w}_t)$

Equivalently  $w_{t+1,j} \leftarrow w_{t,j} - \eta \frac{\partial}{\partial w_{t,j}} f(\mathbf{w}_t) \quad \dots \text{ where } w_j \text{ is the } j_{\text{th}} \text{ weight}$

---

**Algorithm 8.1** Stochastic gradient descent (SGD) update

---

**Require:** Learning rate schedule  $\epsilon_1, \epsilon_2, \dots$

**Require:** Initial parameter  $\theta$

$k \leftarrow 1$

**while** stopping criterion not met **do**

Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

Compute gradient estimate:  $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

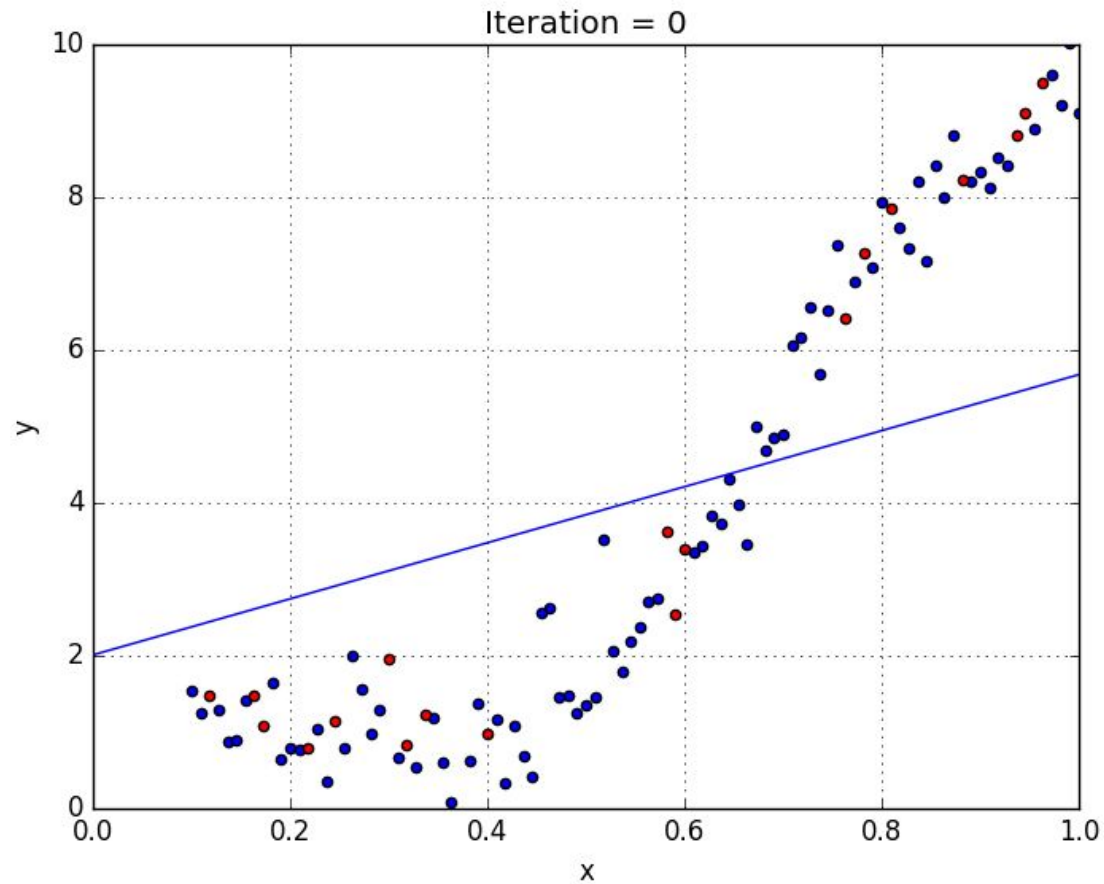
Apply update:  $\theta \leftarrow \theta - \epsilon_k \hat{\mathbf{g}}$

$k \leftarrow k + 1$

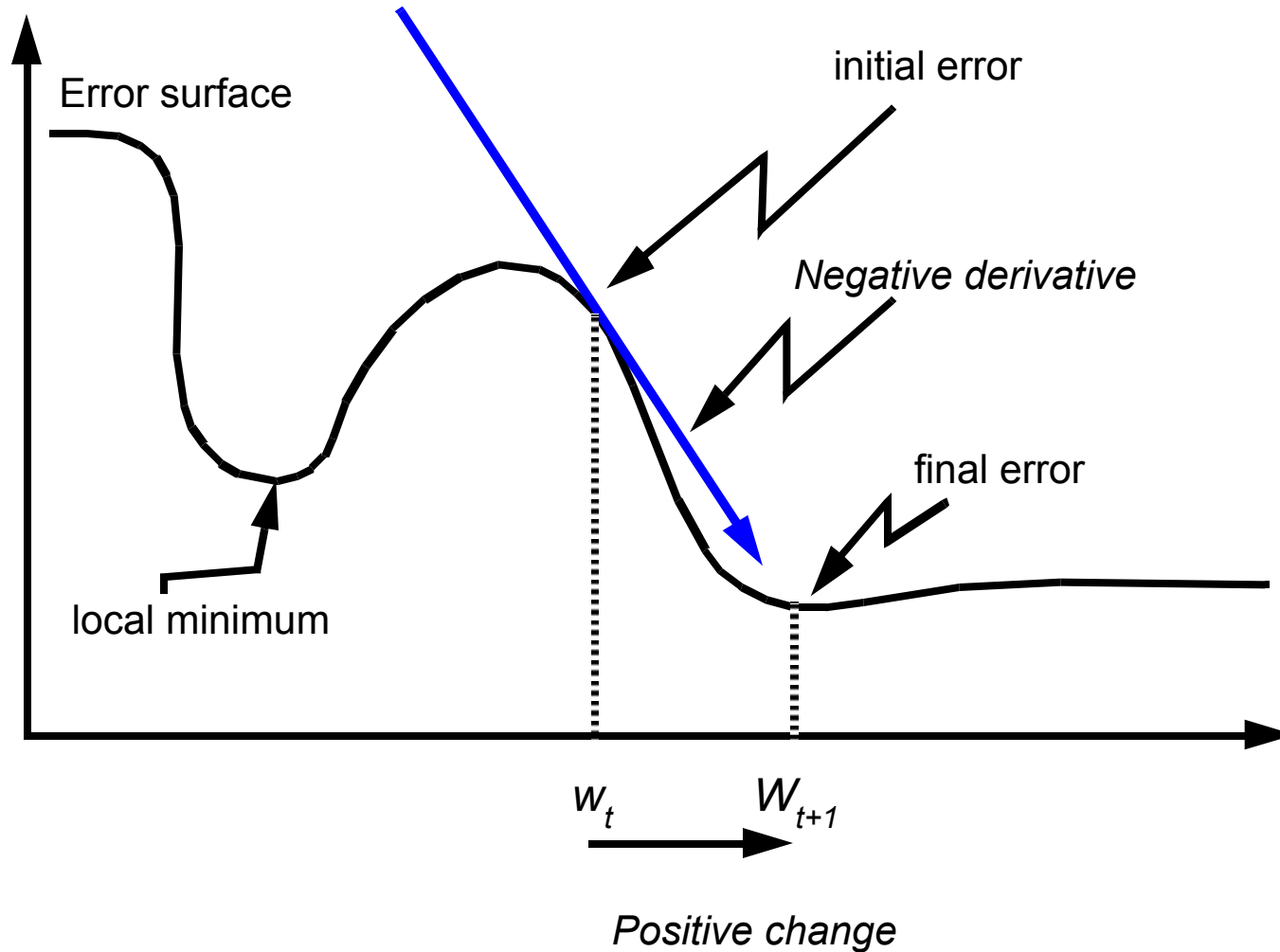
**end while**

---

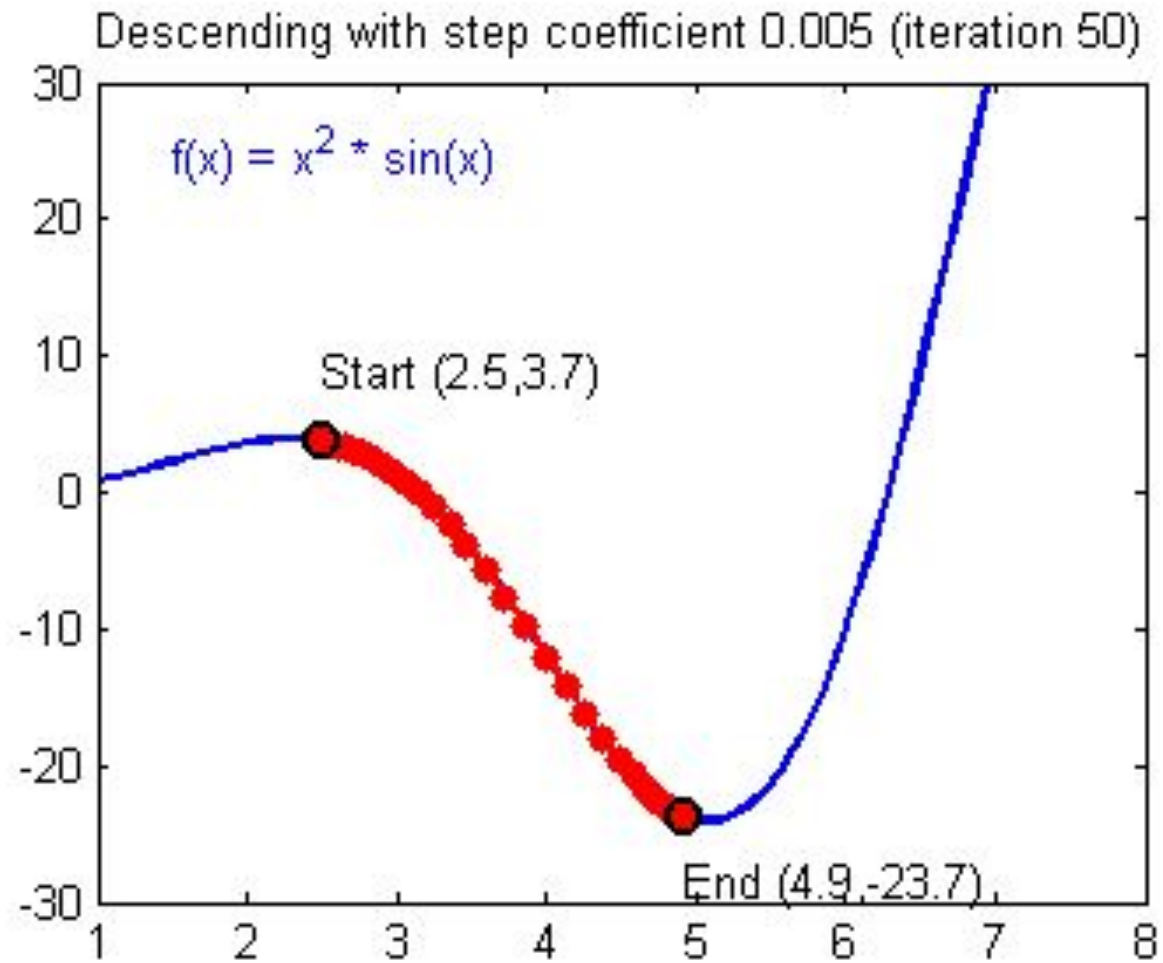
# Using Gradient Descent for Solving Linear Regression



# Minimizing the Error



# Illustration

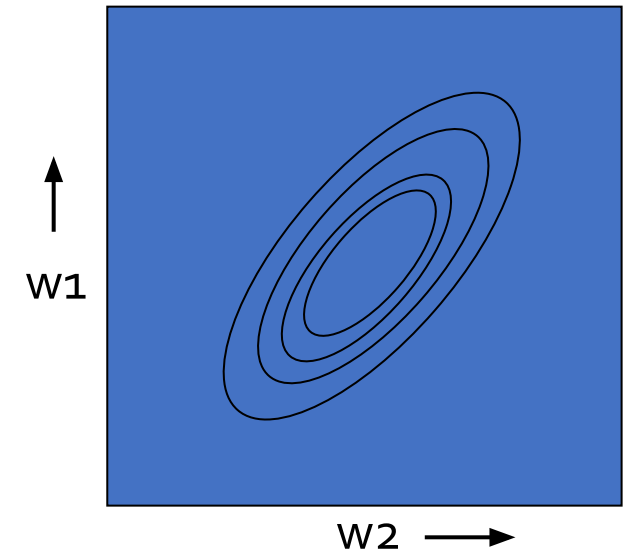
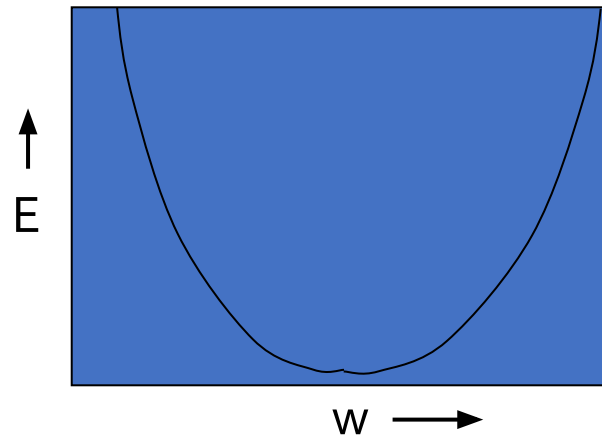
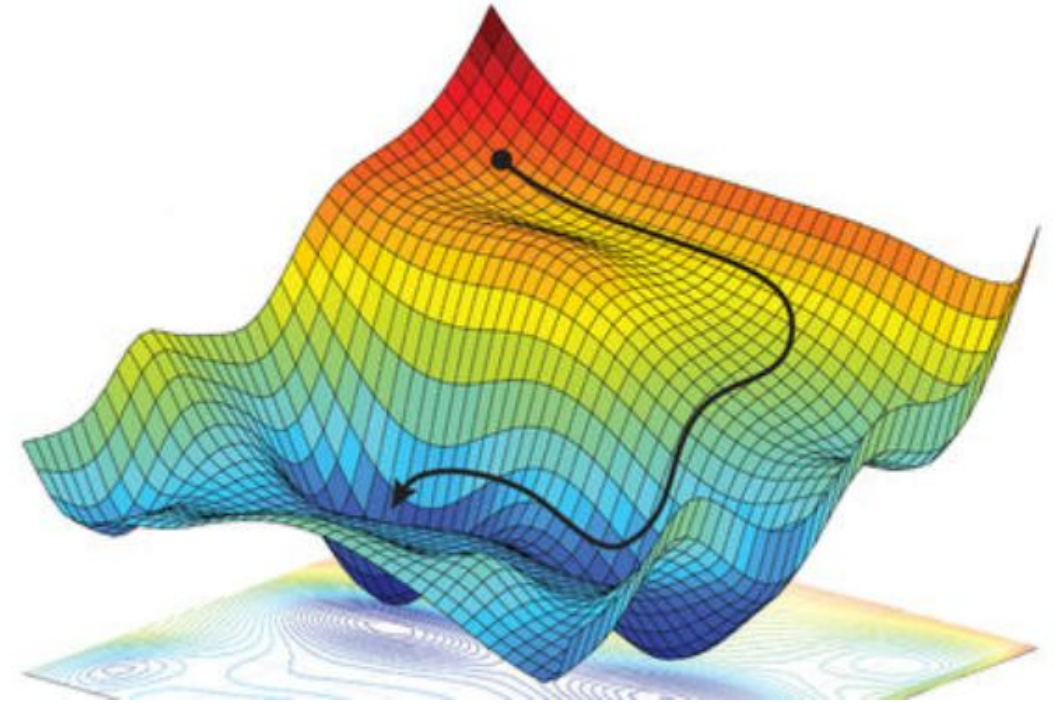


Code in Python can be obtained from: [https://www.cs.toronto.edu/~frossard/post/linear\\_regression/](https://www.cs.toronto.edu/~frossard/post/linear_regression/)



# The error surface

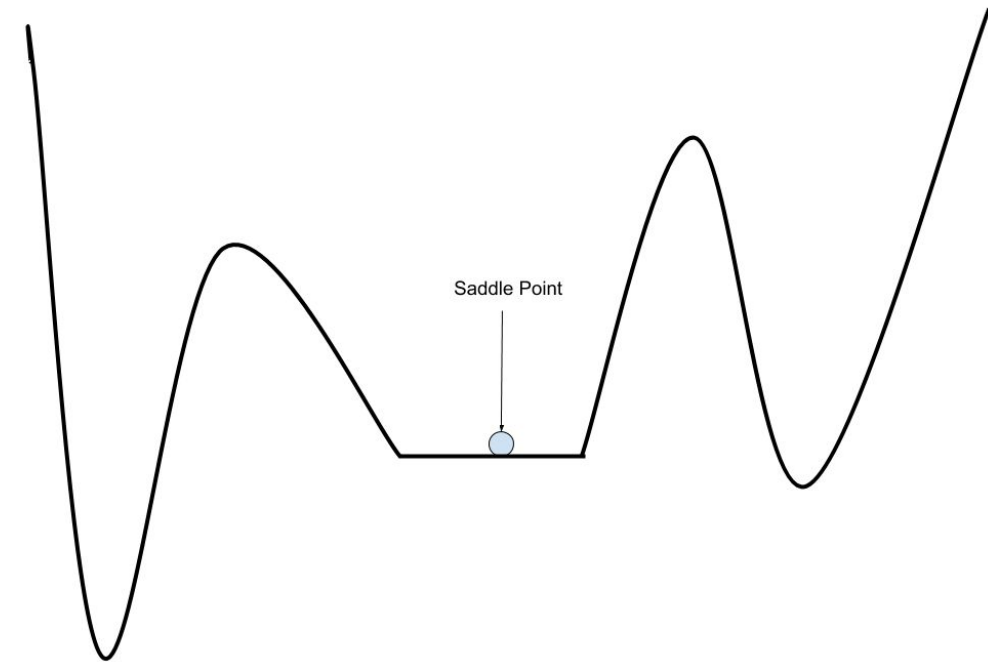
- The error surface lies in a space with a horizontal axis for each weight and one vertical axis for the error.
- It is a quadratic bowl
- Vertical cross-sections are parabolas
- Horizontal cross-sections are ellipses



# The Error Surface

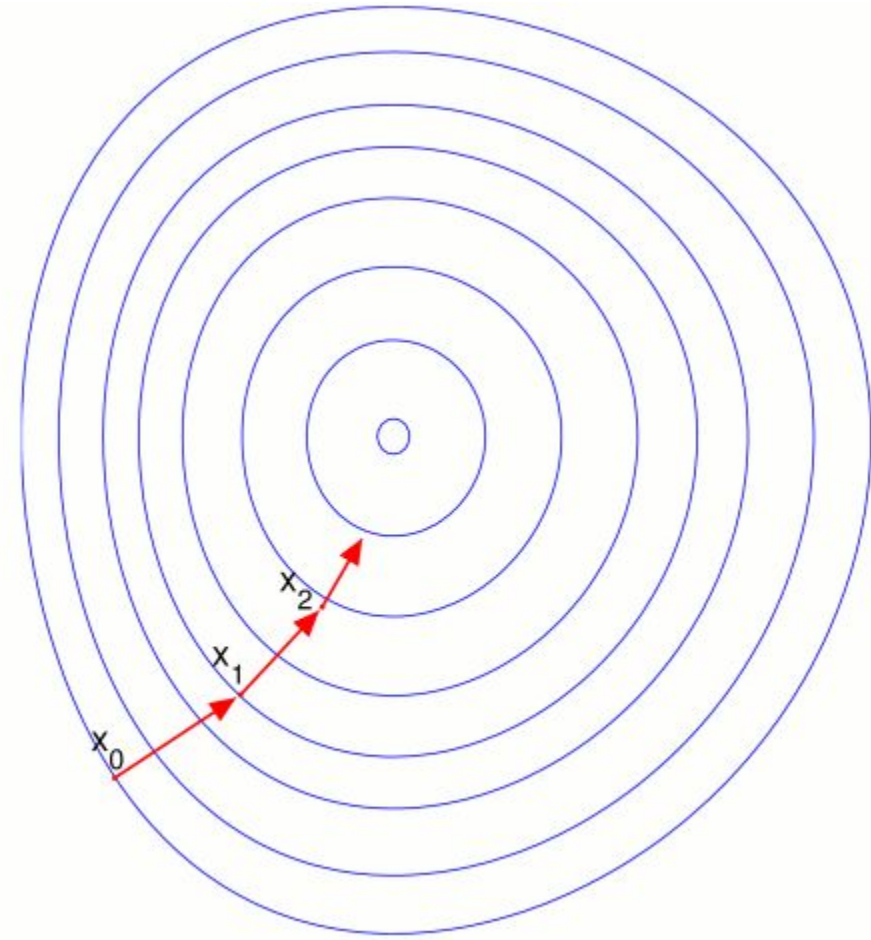
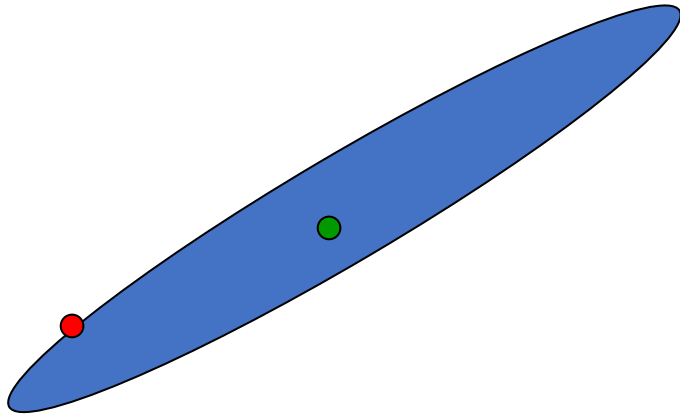
- We only perceive 3 dimensions, but the dimensionality of the error surface is equal to the number of parameters
- Using gradient descent methods mean we might converge to local minima
- Experience shows that in large networks, there are multiple local minima
- BUT, while small networks show a large variance in the loss of the minima, deep networks show a much lower variance (ie there are several possible solutions, but they are all roughly equal)

Hinton: “To deal with hyper-planes in a 14-dimensional space, visualize a 3-D space and say “fourteen” to yourself very loudly. Everyone does it.”



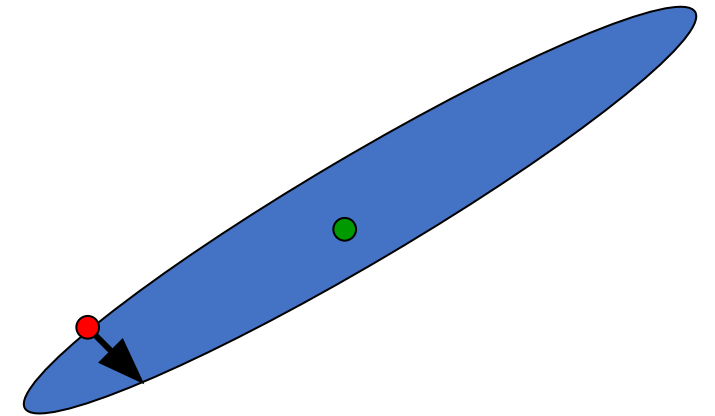
# Illustration of Gradient Descent

- How would gradient descent behave if the error surface is an elongated ellipse?



# Convergence Speed

- The direction of steepest descent does not point at the minimum unless the ellipse is a circle
  - The gradient is large in the direction in which we only want to travel a small distance
  - The gradient is small in the direction in which we want to travel a large distance
- Scaling variables helps turn ellipses into circles

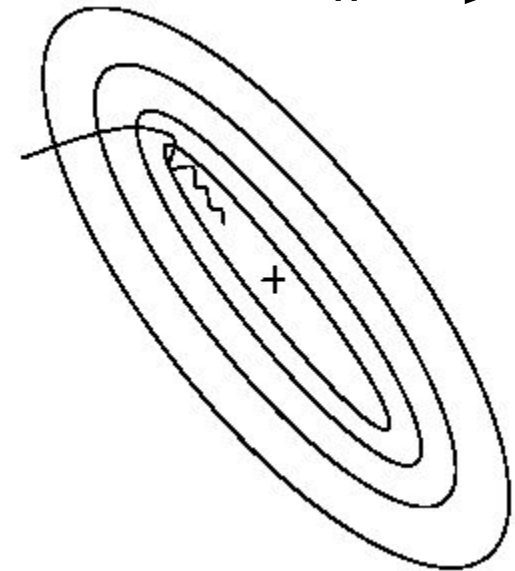
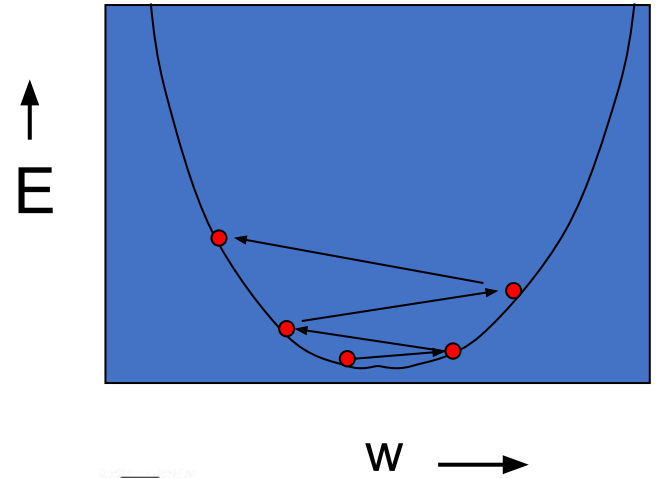
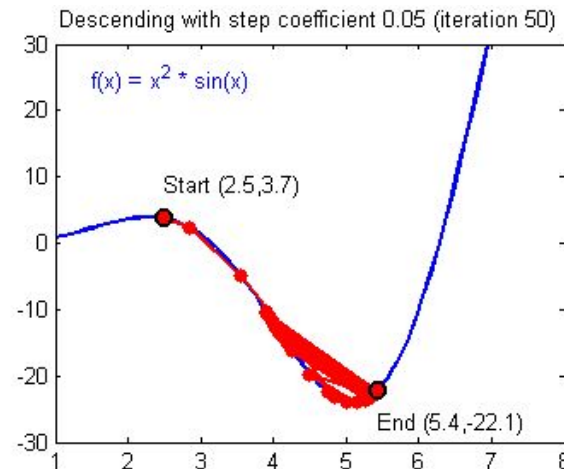
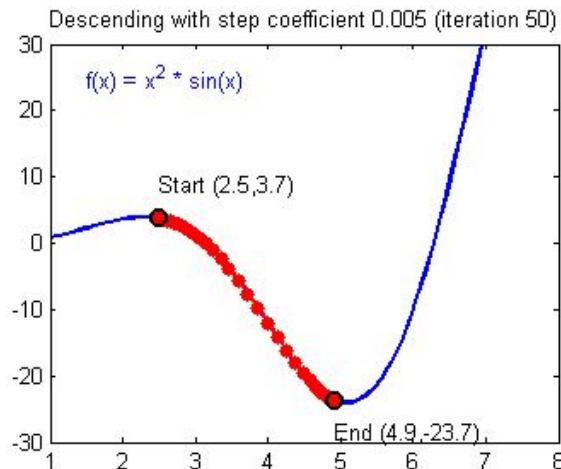


# How learning can go wrong

If the learning rate is too big, it sloshes to and from across the ravine.

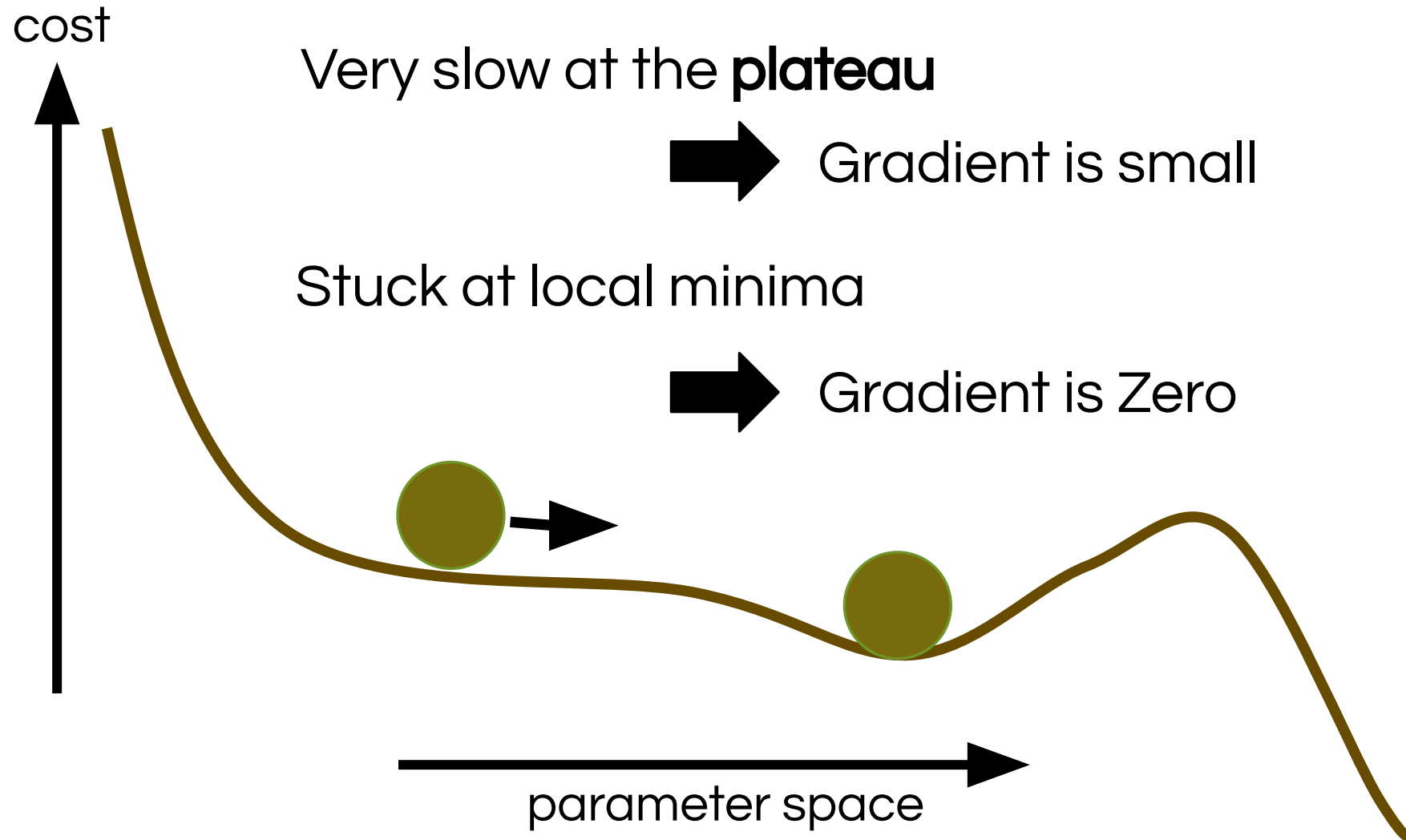
Solution:

- Use small learning rate
- But: if learning rate is too small, it may be slow to coverage



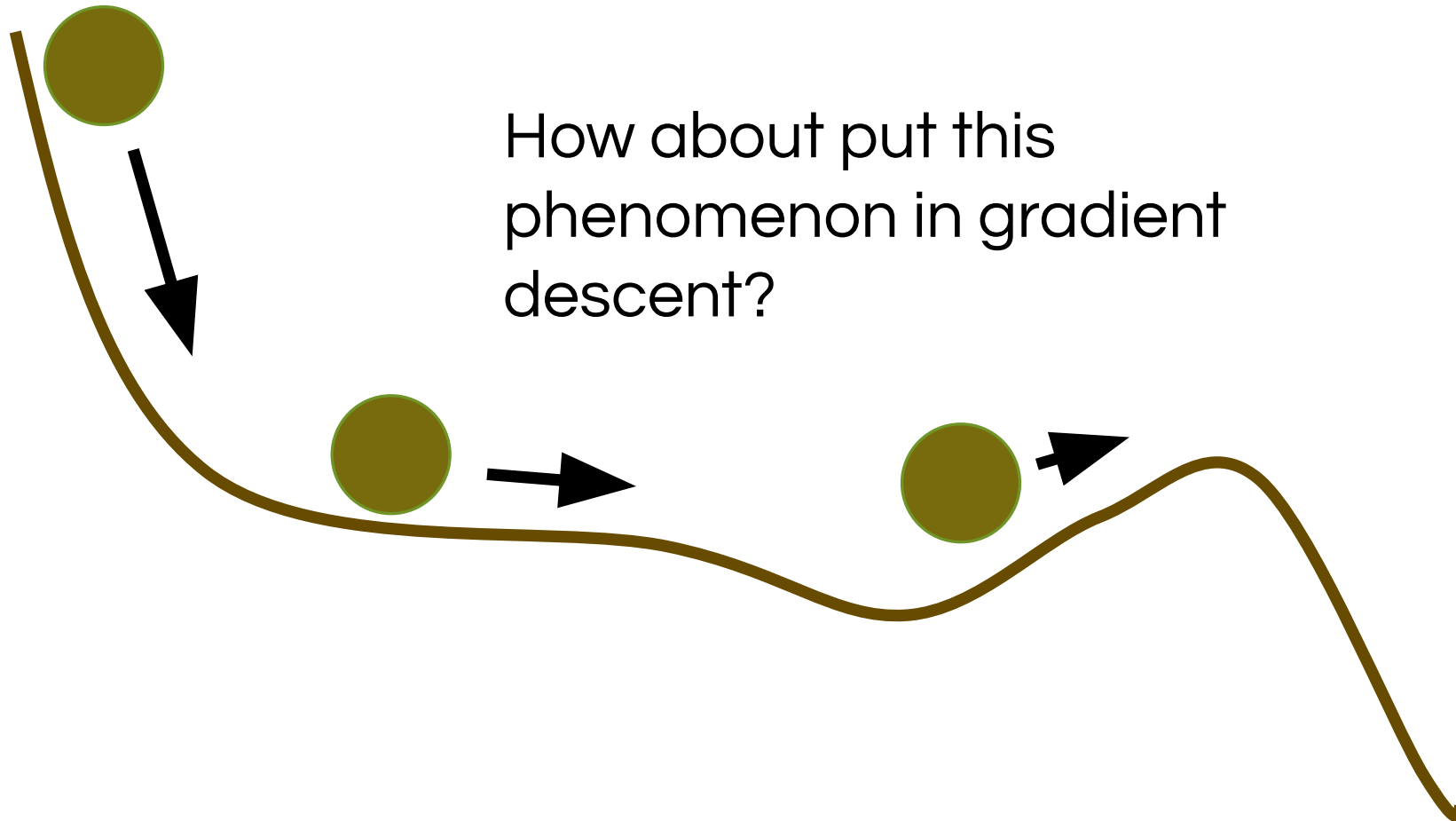
How can we improve  
convergence?

# The problem of gradient descent



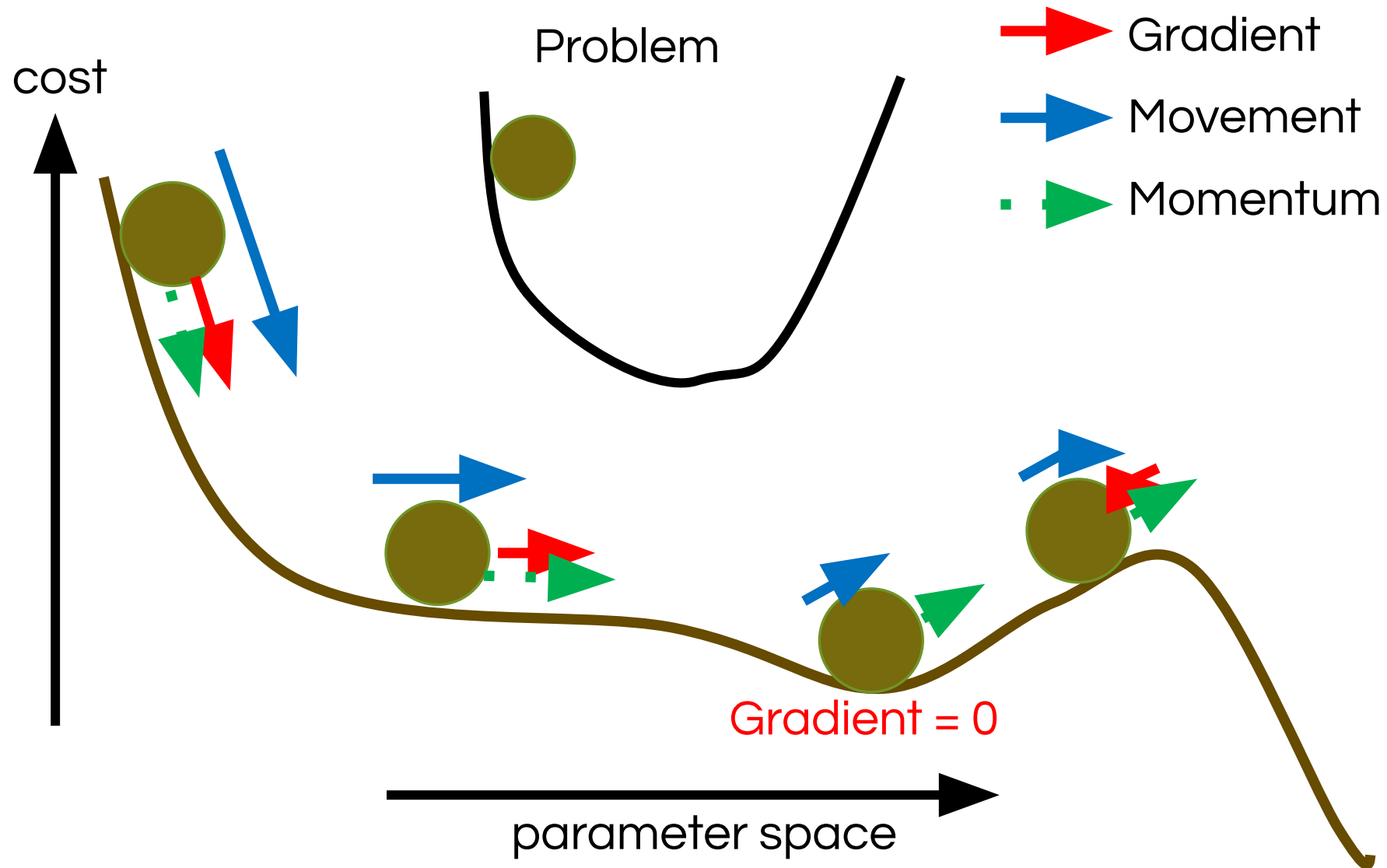
# In the physical world .....

- Momentum

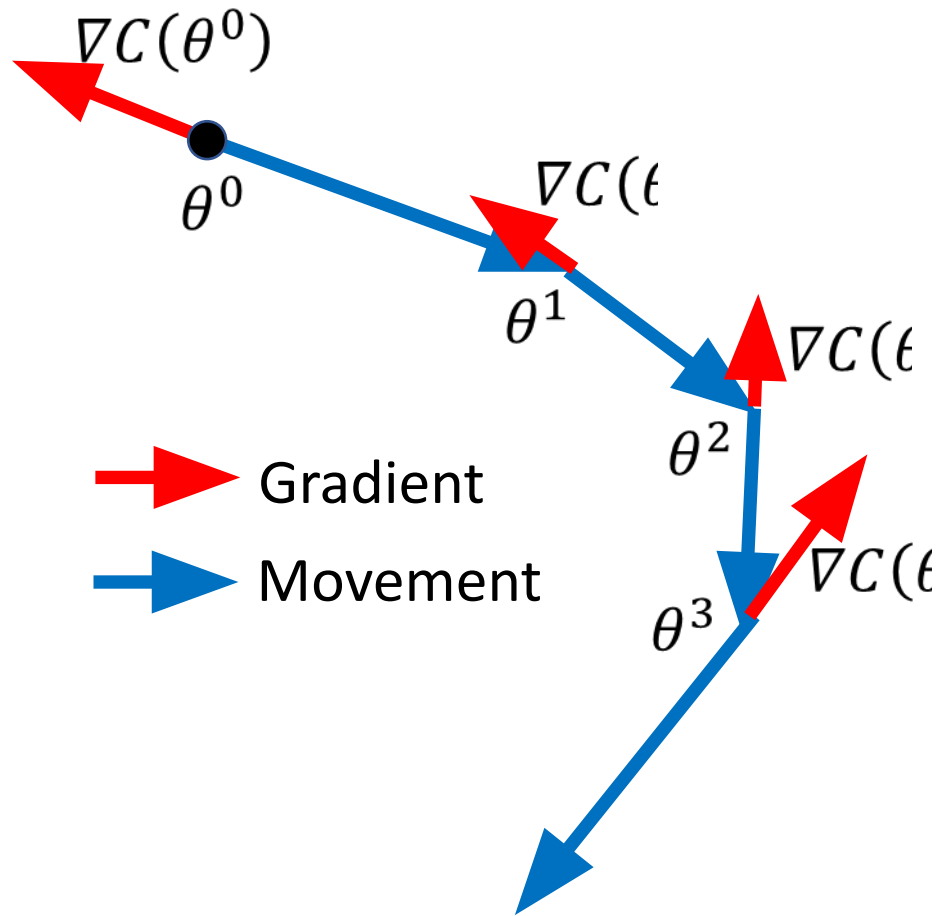




# Gradient descent with Momentum



# Original Gradient descent



Start at position  $\theta^0$

Compute gradient at  $\theta^0$

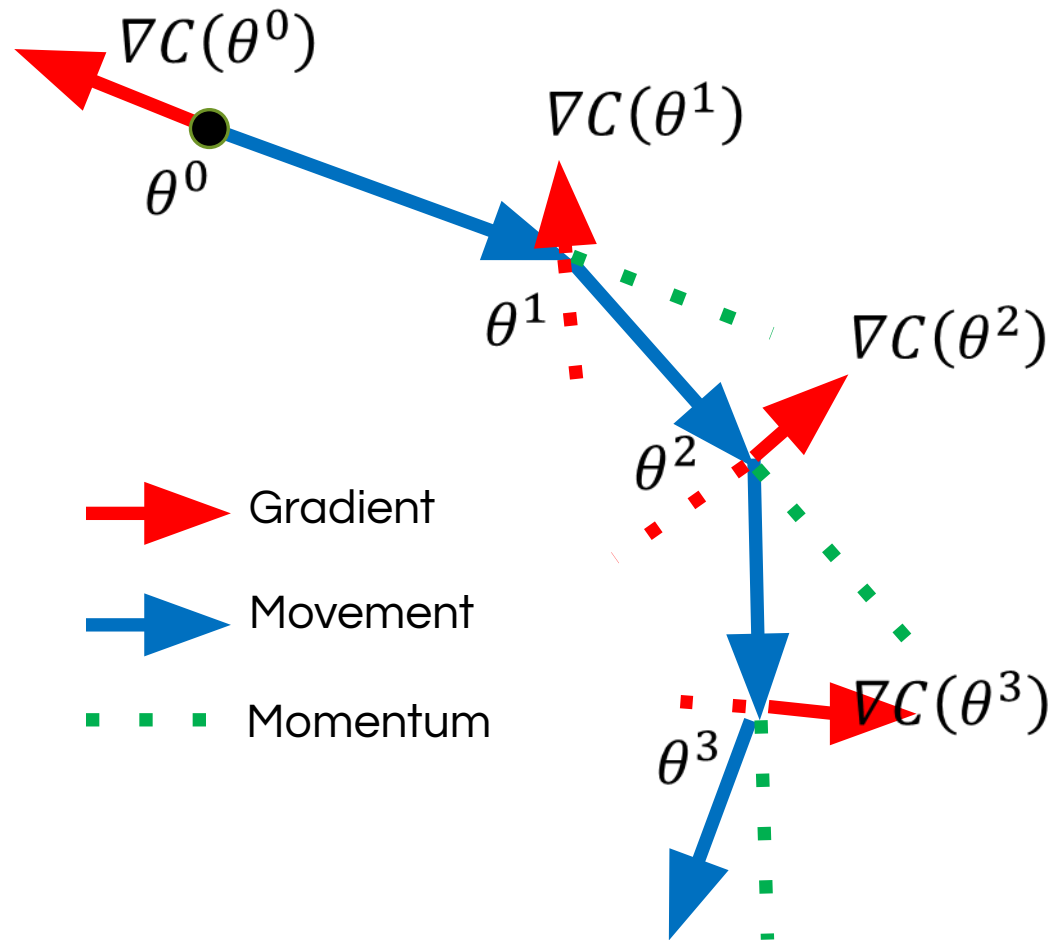
Move to  $\theta^1 = \theta^0 - \mu \nabla C(\theta^0)$

Compute gradient at  $\theta^1$

Move to  $\theta^2 = \theta^1 - \mu \nabla C(\theta^1)$

$\vdots$   
 $\vdots$

# Gradient descent with Momentum (Heavy ball acceleration)



Start at point  $\theta^0$

Momentum  $v^0=0$

( $v^i$ : movement at the  $i$ -th update)

Compute gradient at  $\theta^0$

Momentum  $v^1 = \lambda v^0 - \mu \nabla C(\theta^0)$

Move to  $\theta^1 = \theta^0 + v^1$

Compute gradient at  $\theta^1$

Momentum  $v^2 = \lambda v^1 - \mu \nabla C(\theta^1)$

Move to  $\theta^2 = \theta^1 + v^2$

.....

# Gradient descent with Momentum

$v^i$  is actually the weighted sum of all the previous gradient:

$$\nabla C(\theta^0), \nabla C(\theta^1), \dots \nabla C(\theta^{i-1})$$

$$v^0 = 0$$

$$v^1 = -\mu \nabla C(\theta^0)$$

$$v^2 = -\lambda \mu \nabla C(\theta^0) - \mu \nabla C(\theta^1)$$

.....

Start at point  $\theta^0$

Momentum  $v^0 = 0$

( $v^i$ : movement at the  $i$ -th update)

Compute gradient at  $\theta^0$

$$\text{Momentum } v^1 = \lambda v^0 - \mu \nabla C(\theta^0)$$

$$\text{Move to } \theta^1 = \theta^0 + v^1$$

Compute gradient at  $\theta^1$

$$\text{Momentum } v^2 = \lambda v^1 - \mu \nabla C(\theta^1)$$

$$\text{Move to } \theta^2 = \theta^1 + v^2$$

...

...

# SGD with Momentum

↵

---

**Algorithm 8.2** Stochastic gradient descent (SGD) with momentum

---

**Require:** Learning rate  $\epsilon$ , momentum parameter  $\alpha$

**Require:** Initial parameter  $\theta$ , initial velocity  $v$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient estimate:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .

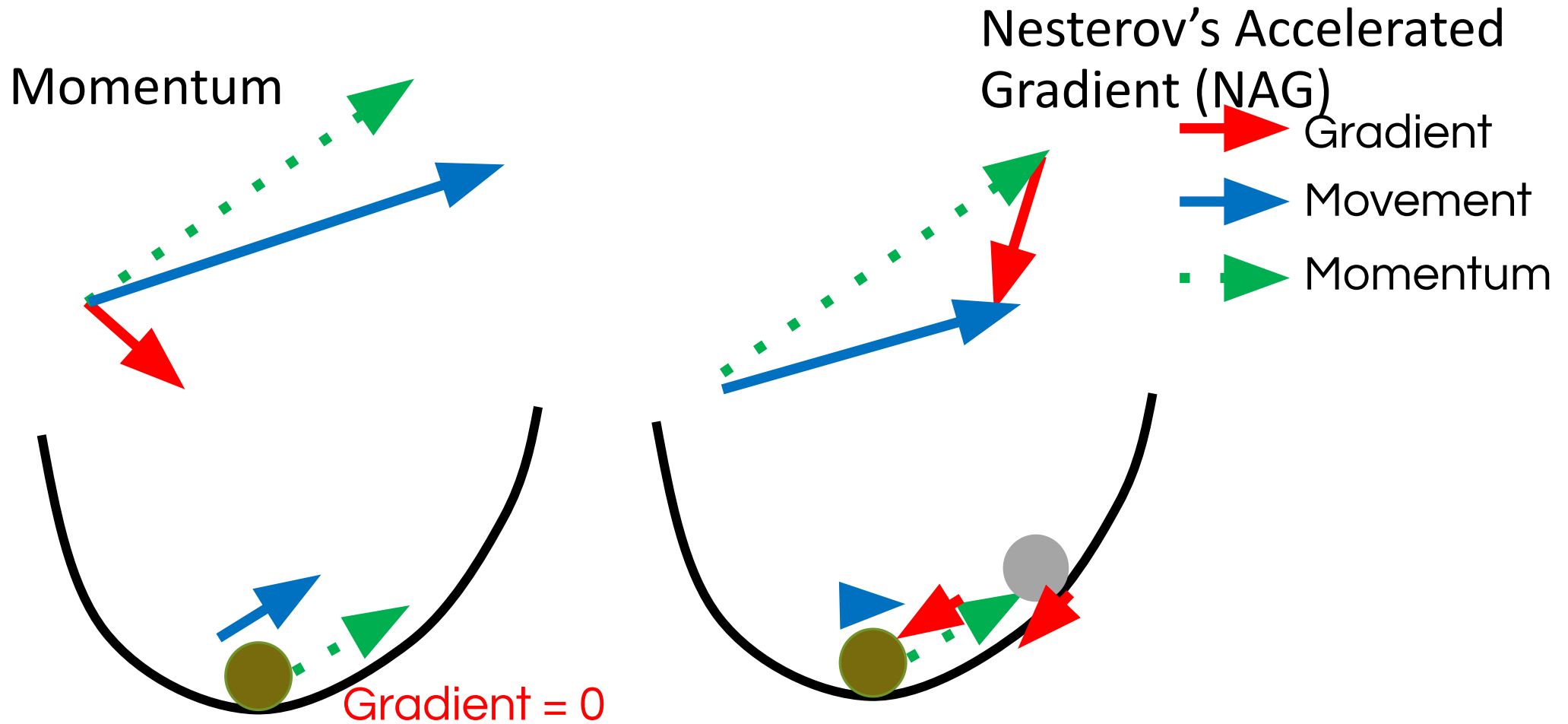
    Compute velocity update:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$ .

    Apply update:  $\theta \leftarrow \theta + \mathbf{v}$ .

**end while**

---

# Nesterov's Accelerated Gradient (NAG)



# SGD with Nesterov Momentum

---

**Algorithm 8.3** Stochastic gradient descent (SGD) with Nesterov momentum

---

**Require:** Learning rate  $\epsilon$ , momentum parameter  $\alpha$

**Require:** Initial parameter  $\theta$ , initial velocity  $v$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding labels  $\mathbf{y}^{(i)}$ .

    Apply interim update:  $\tilde{\theta} \leftarrow \theta + \alpha v$ .

    Compute gradient (at interim point):  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \tilde{\theta}), \mathbf{y}^{(i)})$ .

    Compute velocity update:  $v \leftarrow \alpha v - \epsilon \mathbf{g}$ .

    Apply update:  $\theta \leftarrow \theta + v$ .

**end while**

---

# AdaGrad (Adaptive Gradient)

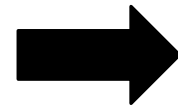
Divide the learning rate by the “*average*” *gradient*

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sigma} g^t$$

$\sigma$ : Average gradient  
of parameter  $w$

Estimated while updating  
the parameters

If  $w$  has small average gradient  $\sigma$



Larger learning rate

If  $w$  has large average gradient  $\sigma$



Smaller learning rate



---

**Algorithm 8.4** The AdaGrad algorithm

---

**Require:** Global learning rate  $\epsilon$

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , perhaps  $10^{-7}$ , for numerical stability

Initialize gradient accumulation variable  $\mathbf{r} = \mathbf{0}$

**while** stopping criterion not met **do**

Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .

Accumulate squared gradient:  $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$ .

Compute update:  $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$ . (Division and square root applied element-wise)

Apply update:  $\theta \leftarrow \theta + \Delta \theta$ .

**end while**

---

# AdaGrad Benefits

- Improves convergence for sparse data (e.g. NLP, Image recognition)
- The accumulated gradient norm increases at every iteration, and so the learning rate shrinks
- AdaGrad eliminates the need to update the learning rate – most often a value of 0.01 is used

# RMSProp (Root Mean Square Propagation)

RMSprop is an extension of Adagrad that deals with its radically diminishing learning rates

Algorithm:

$$\text{Step 1. } G^{(t+1)} = \gamma G^{(t)} + (1 - \gamma) \cdot (g^t)^2$$

*$G$  is a moving exponential average of the square gradient.*

$$\text{Step 2. } w^{(t+1)} = w^{(t)} - \frac{\eta}{\sqrt{G^{(t+1)} + 10^{-7}}} \cdot g^t$$

*Scale the learning rate by  $1/\sqrt{G}$  and follow the negative gradient direction.*

Default settings:  $\gamma = 0.9, \eta = 0.001$ .

# RMSProp

---

**Algorithm 8.5** The RMSProp algorithm

---

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , usually  $10^{-6}$ , used to stabilize division by small numbers

Initialize accumulation variables  $\mathbf{r} = 0$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$ .

    Compute parameter update:  $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$ . ( $\frac{1}{\sqrt{\delta + \mathbf{r}}}$  applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$ .

**end while**

---



# RMSProp with Momentum

---

**Algorithm 8.6** RMSProp algorithm with Nesterov momentum

---

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$ , momentum coefficient  $\alpha$

**Require:** Initial parameter  $\theta$ , initial velocity  $v$

Initialize accumulation variable  $r = 0$

**while** stopping criterion not met **do**

Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

Compute interim update:  $\tilde{\theta} \leftarrow \theta + \alpha v$ .

Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$ .

Accumulate gradient:  $r \leftarrow \rho r + (1 - \rho) g \odot g$ .

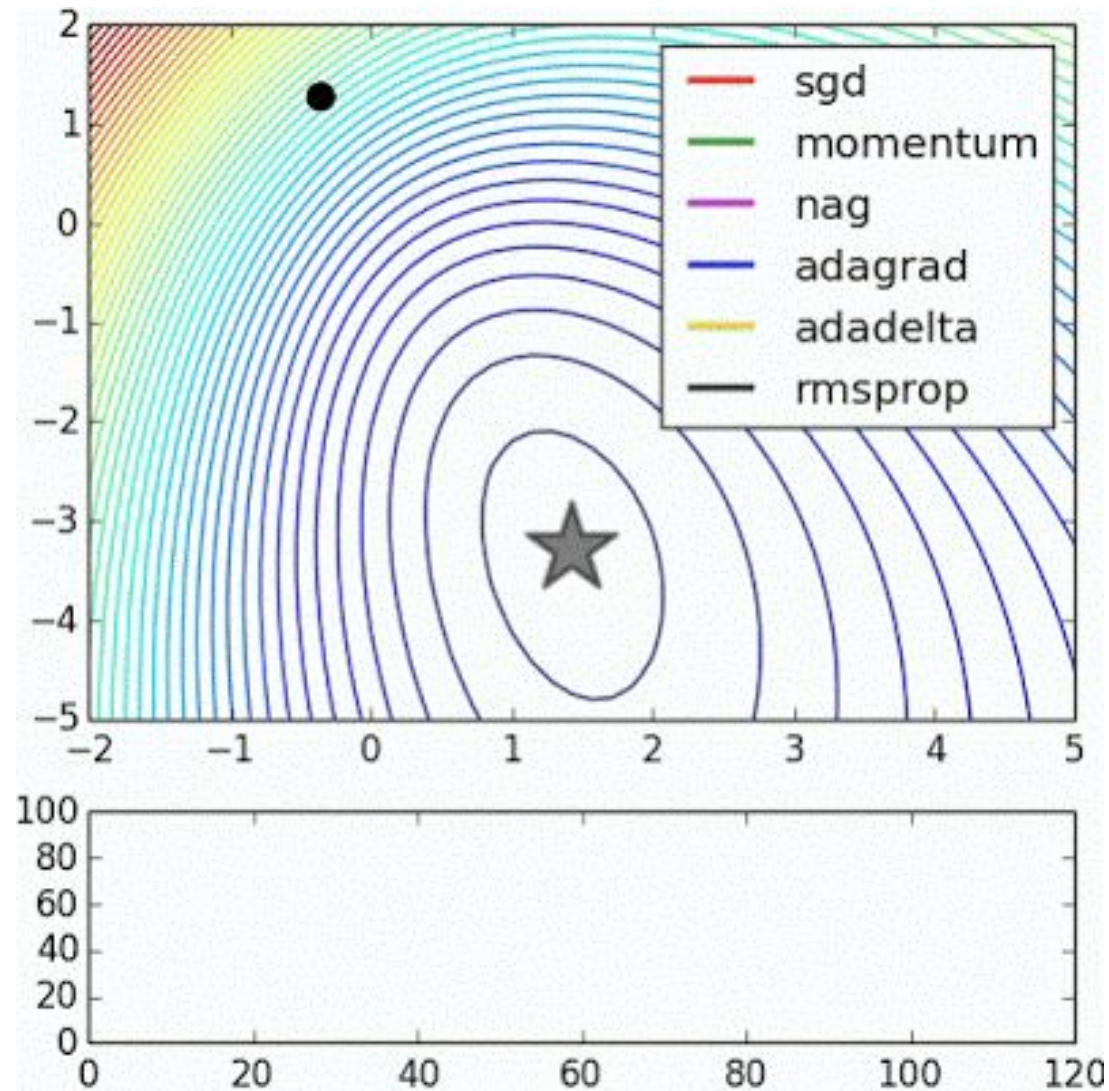
Compute velocity update:  $v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{r}} \odot g$ . ( $\frac{1}{\sqrt{r}}$  applied element-wise)

Apply update:  $\theta \leftarrow \theta + v$ .

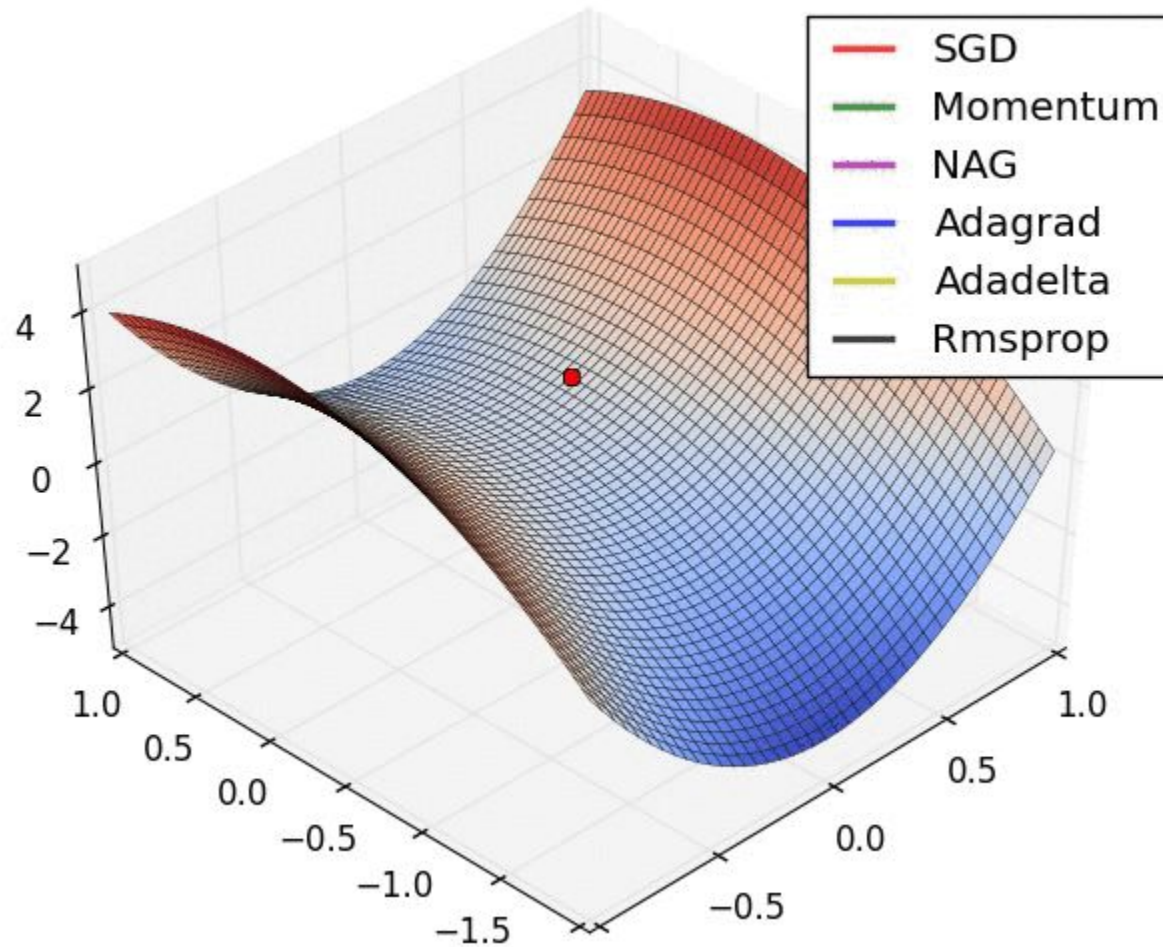
**end while**

---

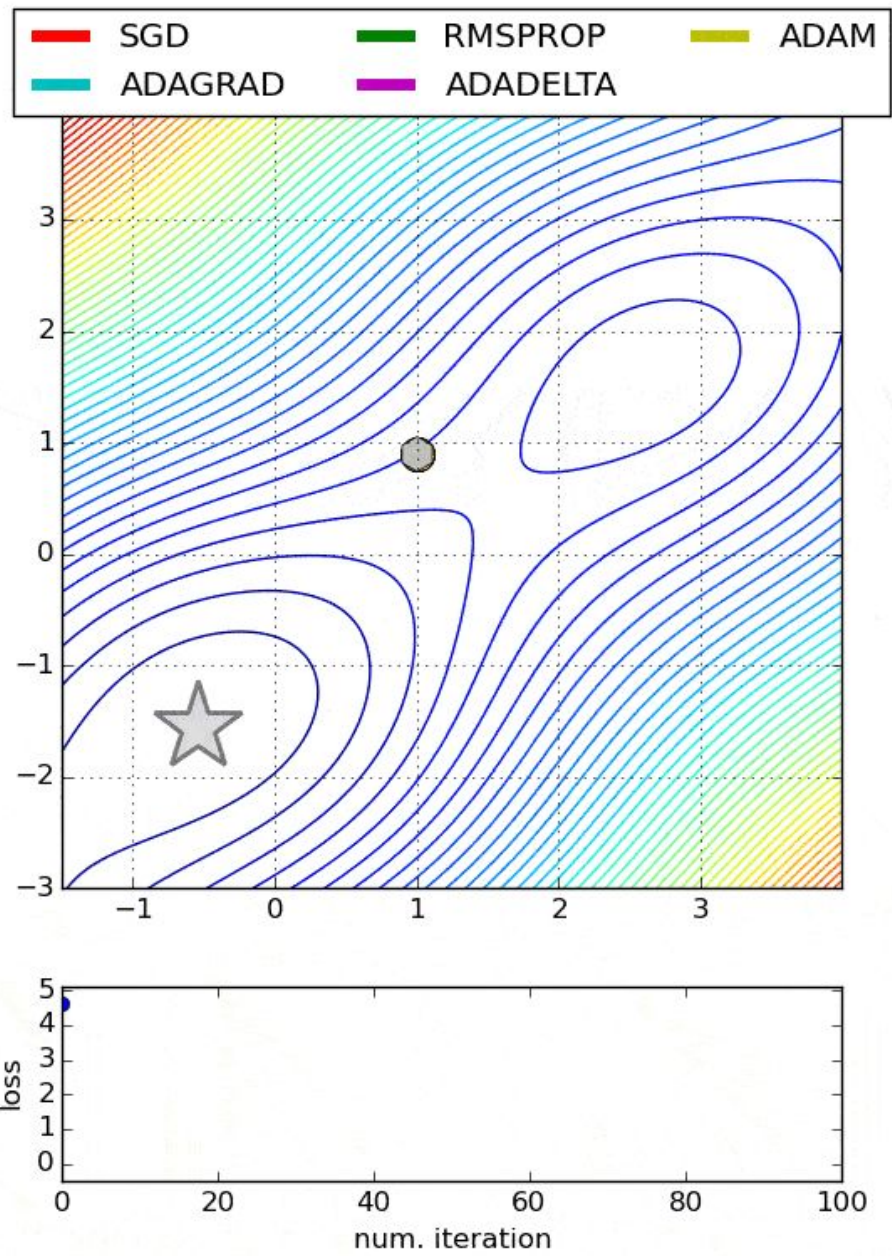
# Many Flavors of Gradient Descent



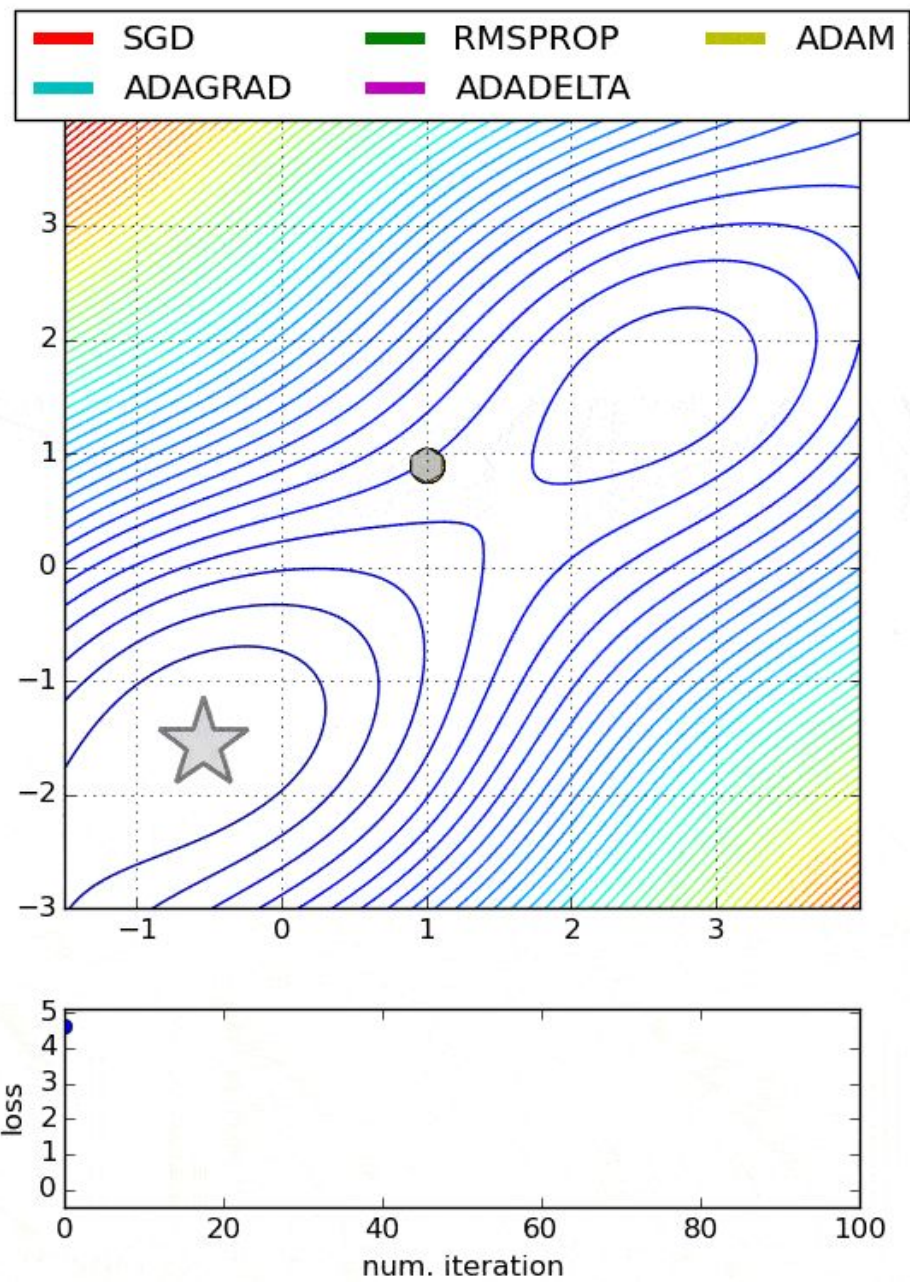
# Many Flavors of Gradient Descent











# Adam (Adaptive Moment Estimation)

- Adaptive learning rate for each parameter
- Keeps exponentially decaying average of past gradients
- Whereas momentum can be seen as a ball running down a slope, Adam behaves like a heavy ball with friction, which thus prefers flat minima in the error surface

# Adam (Adaptive Moment Estimation)

Algorithm [Simplified Version]:

Step 1.  $m^{(t+1)} = \beta_1 m^{(t)} + (1 - \beta_1) \cdot g^t$

$m$  is the exponential moving average of the gradients (similar as the speed in the momentum method);

Step 2.  $v^{(t+1)} = \beta_2 v^{(t)} + (1 - \beta_2) \cdot (g^t)^2$

$v$  is the exponential moving average of the square gradients.

Step 3.  $w^{(t+1)} = w^{(t)} - \frac{\eta}{\sqrt{v^{(t+1)} + 10^{-8}}} \cdot m^{(t+1)}$

Default settings:  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\eta = 0.001$ .

---

**Algorithm 8.7** The Adam algorithm

---

**Require:** Step size  $\epsilon$  (Suggested default: 0.001)

**Require:** Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1)$ .  
(Suggested defaults: 0.9 and 0.999 respectively)

**Require:** Small constant  $\delta$  used for numerical stabilization (Suggested default:  $10^{-8}$ )

**Require:** Initial parameters  $\theta$

Initialize 1st and 2nd moment variables  $\mathbf{s} = \mathbf{0}$ ,  $\mathbf{r} = \mathbf{0}$

Initialize time step  $t = 0$

**while** stopping criterion not met **do**

Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

Update biased first moment estimate:  $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

Update biased second moment estimate:  $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

Correct bias in first moment:  $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

Correct bias in second moment:  $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

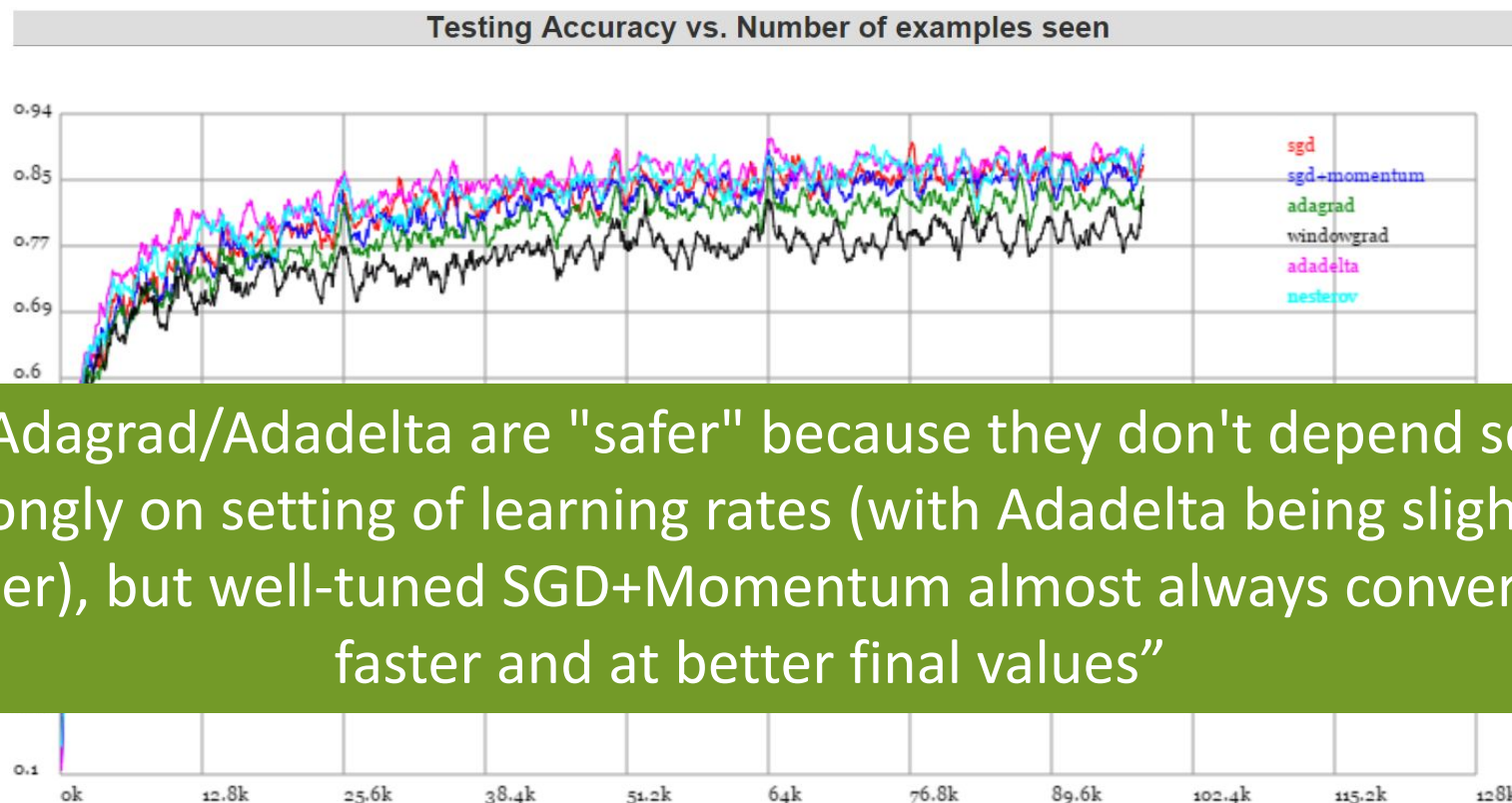
Compute update:  $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$  (operations applied element-wise)

Apply update:  $\theta \leftarrow \theta + \Delta \theta$

**end while**

---

# CNN training with different optimization algorithms



“Adagrad/Adadelta are "safer" because they don't depend so strongly on setting of learning rates (with Adadelta being slightly better), but well-tuned SGD+Momentum almost always converges faster and at better final values”



# Comparing GD Methods

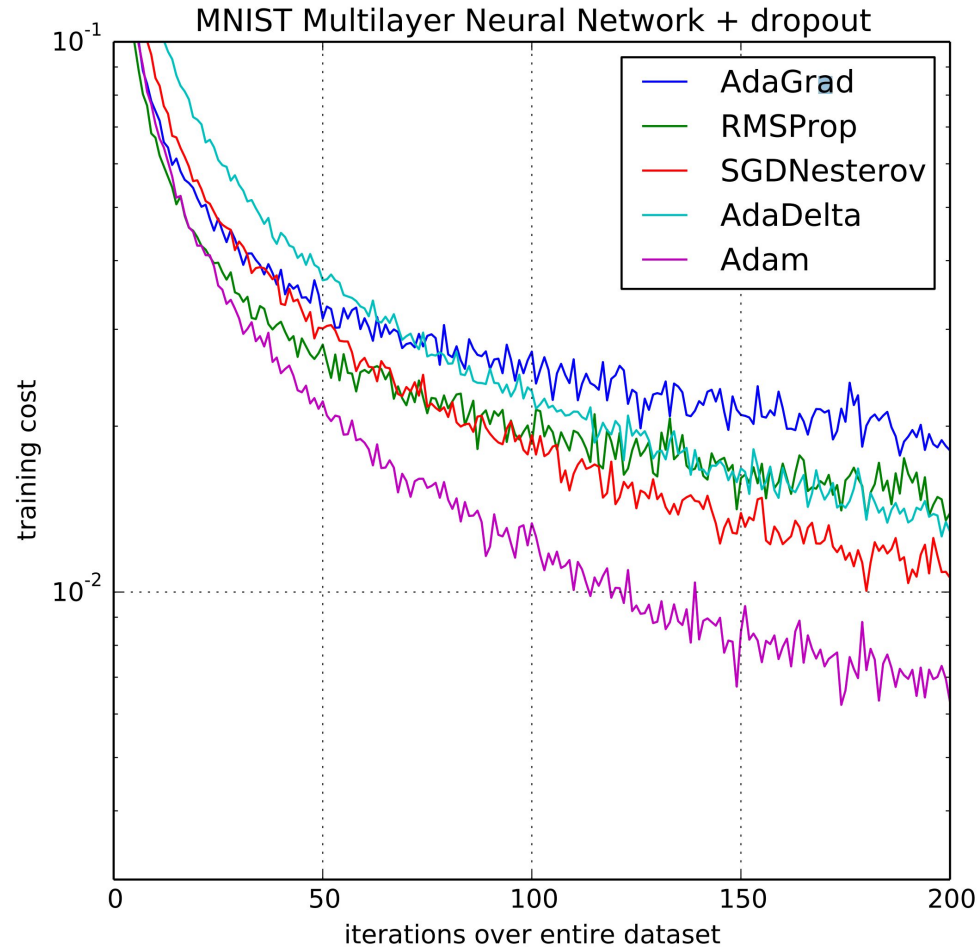


Image credit: [Kingma & Ba, 2015](#)

## Network configurations:

Input layer: 28x28

2 hidden layers (1000 ReLUs each)

Output: 10 SoftMax

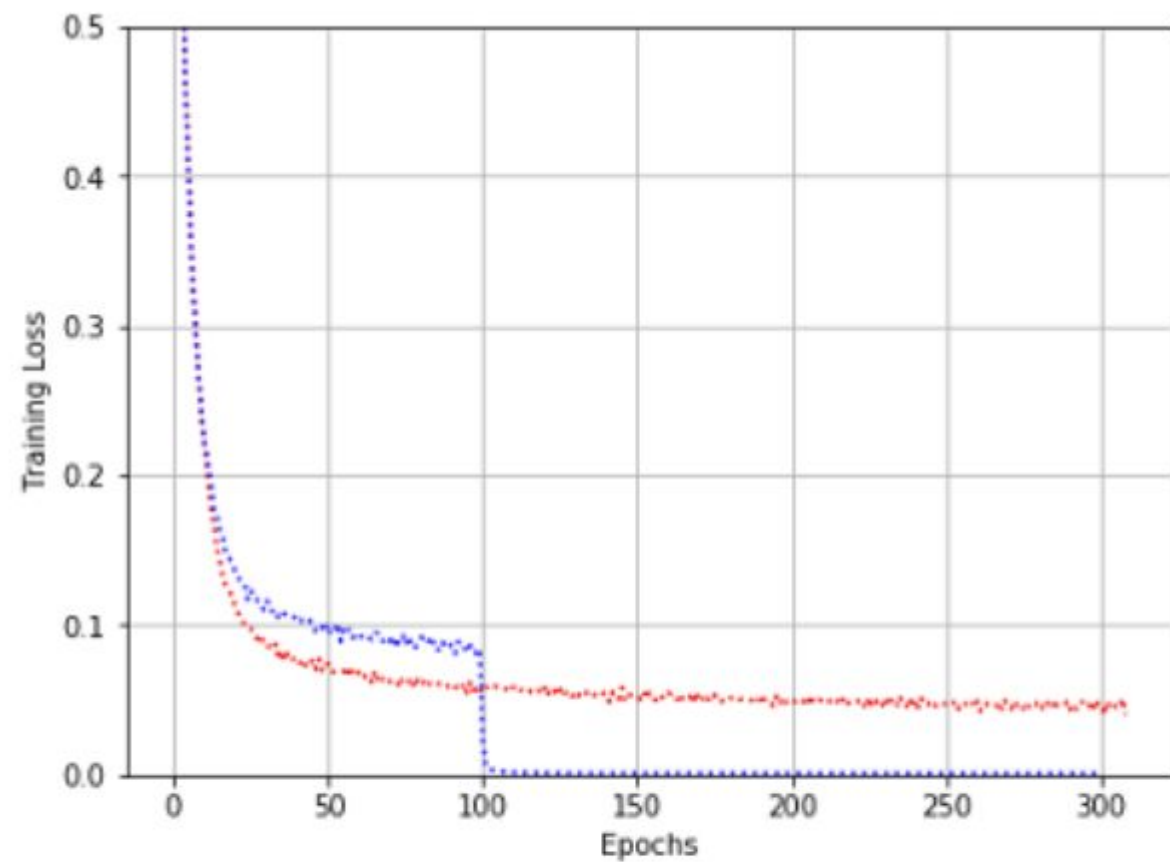
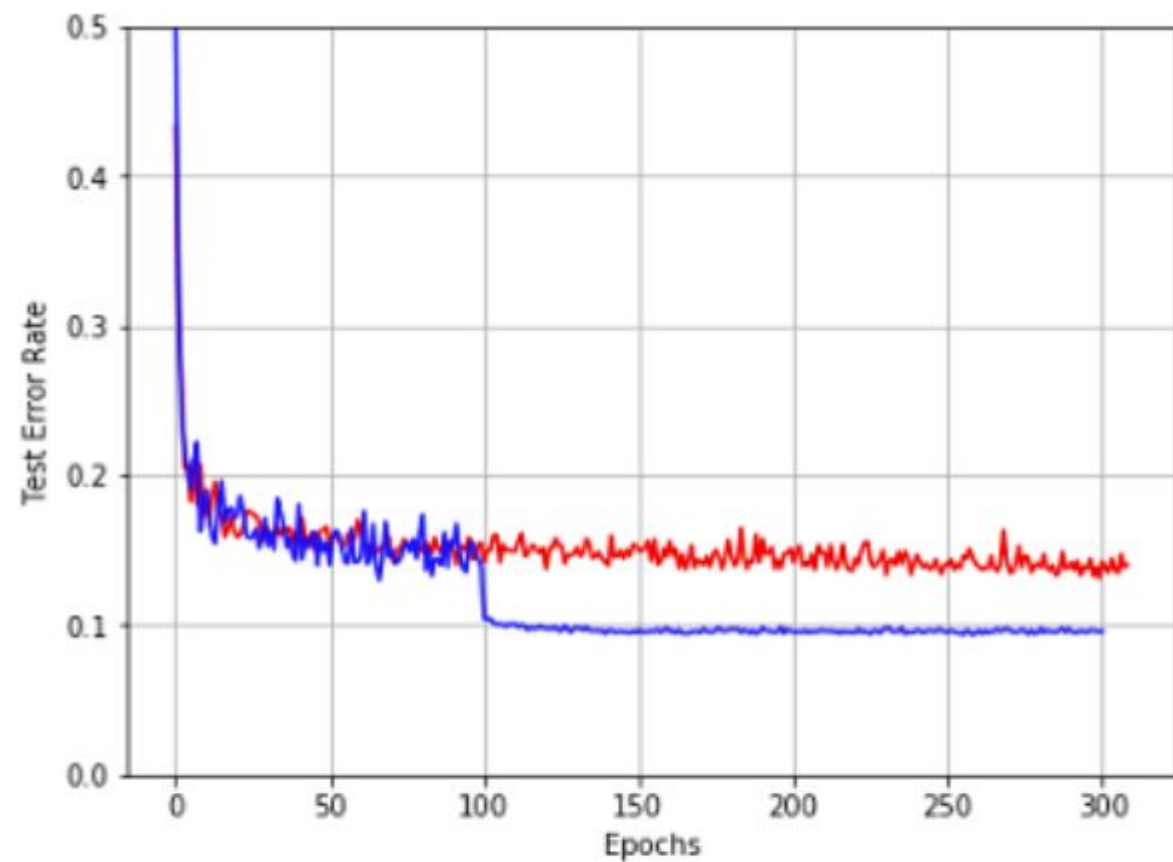
Cross-entropy loss.

In practice, Adam is the default choice to use

[\[Fei-Fei Li, Andrej Karpathy, Justin Johnson\]](#)

# No Free Lunch

- RMSprop is an extension of Adagrad that deals with its radically diminishing learning rates. It is identical to Adadelta, except that Adadelta uses the RMS of parameter updates in the nominator update rule. Adam, finally, adds bias-correction and momentum to RMSprop.
- RMSprop, Adadelta, and Adam are very similar algorithms that do well in similar circumstances
- Adam has become the default algorithm used
- But there are cases in which GD and SGD provide much better results:
  - The number of parameters exceeds the number of data points
  - The rate decay of SGD is well-tuned
- Some modern methods combine Adam for initial stages with SGD for later ones



**Fig. 5:** Comparison of the results obtained by tuning and not tuning the learning rate. **Red lines:** the results without tuning; **Blue lines:** the results with tuning.



# AdamW

- Modifies the typical implementation of weight decay in Adam, by decoupling weight decay from the gradient update

---

## Algorithm 2 Adam with $L_2$ regularization and Adam with weight decay (AdamW)

---

```
1: given  $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, w \in \mathbb{R}$ 
2: initialize time step  $t \leftarrow 0$ , parameter vector  $\mathbf{x}_{t=0} \in \mathbb{R}^n$ , first
   moment vector  $\mathbf{m}_{t=0} \leftarrow \mathbf{0}$ , second moment vector  $\mathbf{v}_{t=0} \leftarrow \mathbf{0}$ ,
   schedule multiplier  $\eta_{t=0} \in \mathbb{R}$ 
3: repeat
4:    $t \leftarrow t + 1$ 
5:    $\nabla f_t(\mathbf{x}_{t-1}) \leftarrow \text{SelectBatch}(\mathbf{x}_{t-1})$   $\triangleright$  select batch and
     return the corresponding gradient
6:    $\mathbf{g}_t \leftarrow \nabla f_t(\mathbf{x}_{t-1}) + w\mathbf{x}_{t-1}$ 
7:    $\mathbf{m}_t \leftarrow \beta_1\mathbf{m}_{t-1} + (1 - \beta_1)\mathbf{g}_t$   $\triangleright$  here and below all
     operations are element-wise
8:    $\mathbf{v}_t \leftarrow \beta_2\mathbf{v}_{t-1} + (1 - \beta_2)\mathbf{g}_t^2$ 
9:    $\hat{\mathbf{m}}_t \leftarrow \mathbf{m}_t / (1 - \beta_1^t)$   $\triangleright \beta_1$  is taken to the power of  $t$ 
10:   $\hat{\mathbf{v}}_t \leftarrow \mathbf{v}_t / (1 - \beta_2^t)$   $\triangleright \beta_2$  is taken to the power of  $t$ 
11:   $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$   $\triangleright$  can be fixed, decay, or
     also be used for warm restarts
12:   $\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \eta_t \left( \alpha \hat{\mathbf{m}}_t / (\sqrt{\hat{\mathbf{v}}_t} + \epsilon) + w\mathbf{x}_{t-1} \right)$ 
13: until stopping criterion is met
14: return optimized parameters  $\mathbf{x}_t$ 
```

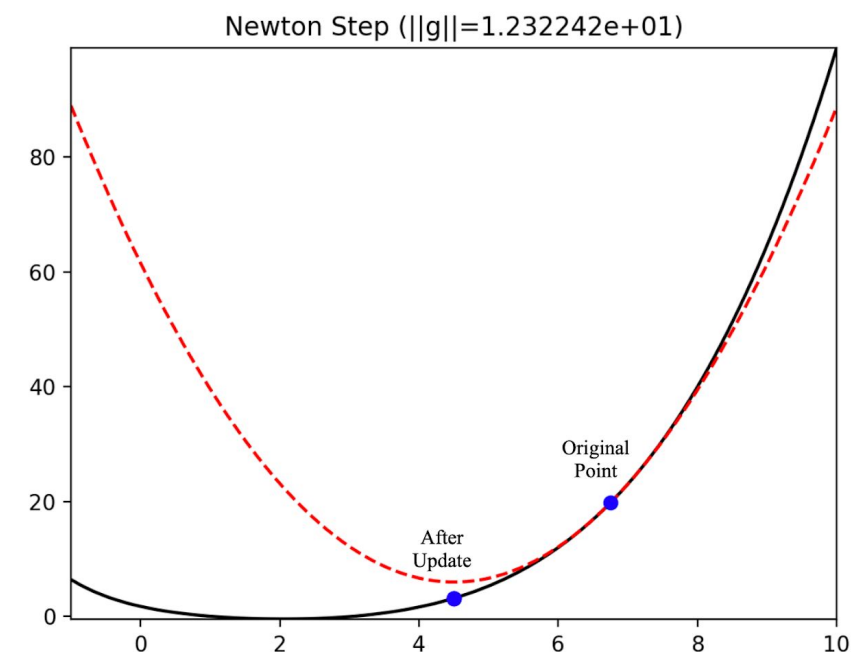
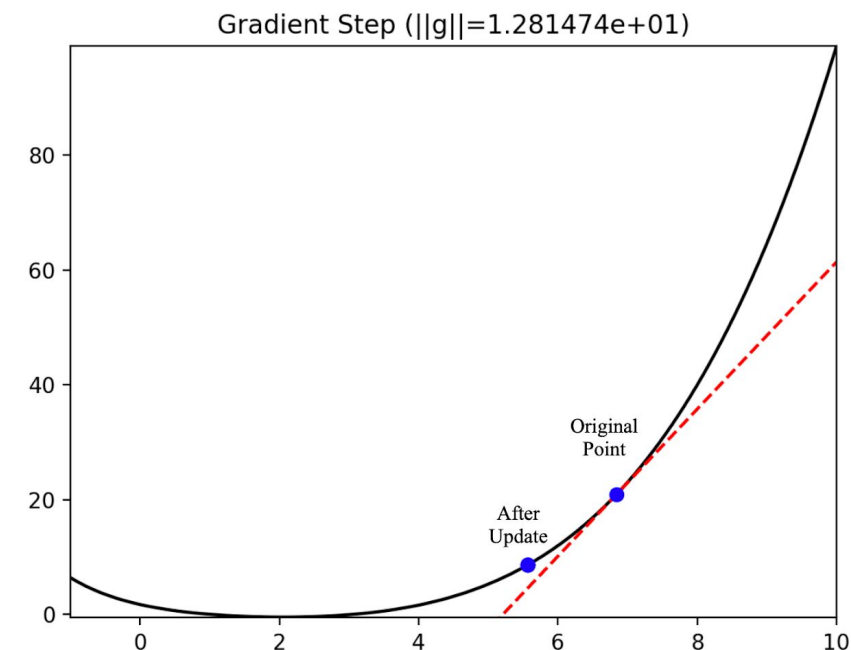
---

# Alternatives to Gradient Descent

- The gradient provides the slope of the function in each dimension
- More information is included in the Hessian – the square matrix of second order derivatives
- This is called *Newton's method*

$$x \leftarrow x - [Hf(x)]^{-1} \nabla f(x)$$

- The Hessian describes the local curvature of the loss function, which allows us to perform a more efficient update
- Leads to more aggressive steps in directions of shallow curvature and shorter steps in directions of steep curvature
- Note that there is no learning rate hyperparameter(s) – an advantage



# Using the Hessian – Newton's Method

- Computing the Hessian requires computing many parameters – assuming 1M parameters, the Hessian has  $1M \times 1M$  values
- Current research focuses on approximating the inverse Hessian, preferably on mini-batches and not the entire dataset
- E.g. [L-BFGS](#) uses historical gradient information to approximate the Hessian
- In practice, SGD+Momentum methods are much more common for optimizing neural networks
- The good news – you generally don't have to worry about any of that

# What it means for you?

- A good optimizer can help you squeeze more juice from the lemon
- But too often people start optimizing too early
- Hyper-parameter optimization should be among the last things you do
  - Explore the data and errors
  - Consider getting more data
  - Evaluate your loss function
  - ...

# Neural Networks in PyTorch

- We define the loss function and the optimizer:

```
import torch.optim as optim

criterion = nn.CrossEntropyLoss()
optimizer = optim.SGD(net.parameters(), lr=0.001, momentum=0.9)
```

- PyTorch Optimizers

- `public torch::optim::Adagrad` (Class Adagrad)
- `public torch::optim::Adam` (Class Adam)
- `public torch::optim::AdamW` (Class AdamW)
- `public torch::optim::LBFGS` (Class LBFGS)
- `public torch::optim::RMSprop` (Class RMSprop)
- `public torch::optim::SGD` (Class SGD)

# Training the Network is done in a loop

```
for epoch in range(2): # loop over the dataset multiple times

    running_loss = 0.0
    for i, data in enumerate(trainloader, 0):
        # get the inputs; data is a list of [inputs, labels]
        inputs, labels = data

        # zero the parameter gradients
        optimizer.zero_grad()

        # forward + backward + optimize
        outputs = net(inputs)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()

        # print statistics
        running_loss += loss.item()
        if i % 2000 == 1999: # print every 2000 mini-batches
            print('[%d, %5d] loss: %.3f' %
                  (epoch + 1, i + 1, running_loss / 2000))
            running_loss = 0.0

    print('Finished Training')
```

# Per-layer Optimizer Params

```
optim.SGD([
    {'params': model.base.parameters()},
    {'params': model.classifier.parameters(), 'lr': 1e-3}
], lr=1e-2, momentum=0.9)
```

This means that `model.base`'s parameters will use the default learning rate of `1e-2`, `model.classifier`'s parameters will use a learning rate of `1e-3`, and a momentum of `0.9` will be used for all parameters.



# Suggested Reading

- Chapter 8 of the book “Deep Learning” –  
<https://www.deeplearningbook.org/contents/optimization.html>

# Recommended Resources

- <https://hackernoon.com/some-state-of-the-art-optimizers-in-neural-networks-a3c2ba5a5643>
- <https://runder.io/optimizing-gradient-descent/>
- [AdaBound –](#)  
<https://medium.com/syncedreview/iclr-2019-fast-as-adam-good-as-sgd-new-optimizer-has-both-78e37e8f9a34>