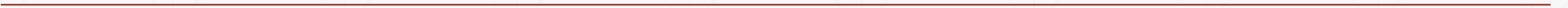# 动量与学习率

主讲：龙良曲

# Outline

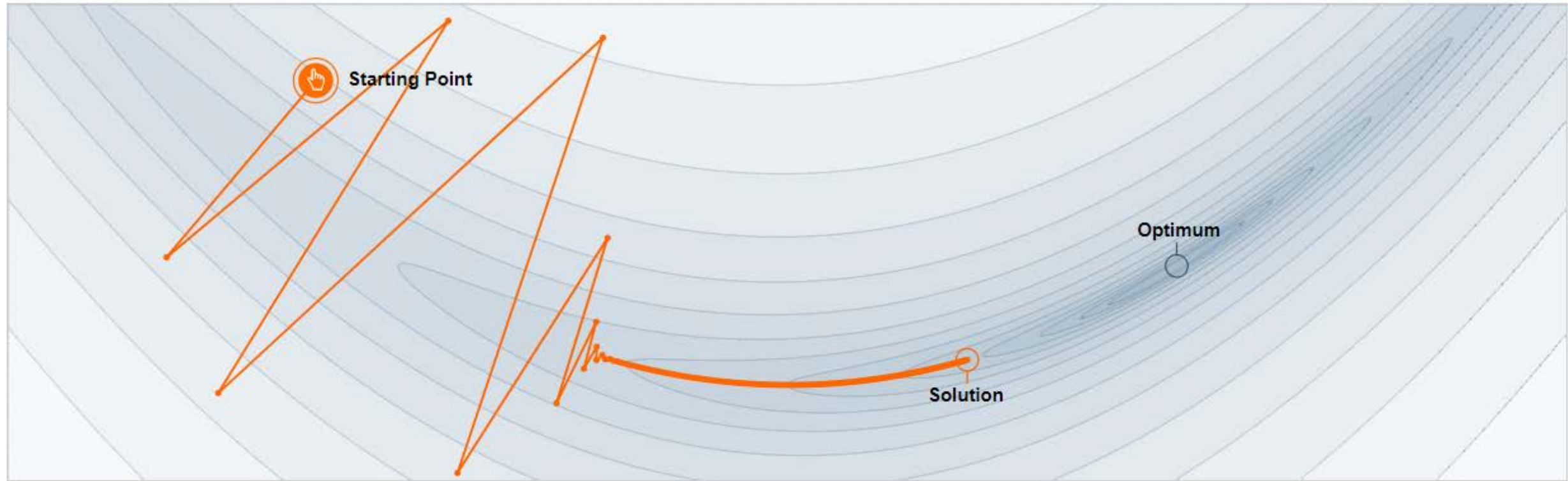- momentum

- learning rate decay

# Momentum

$$w^{k+1} = w^k - \alpha \nabla f(w^k).$$

$$z^{k+1} = \beta z^k + \nabla f(w^k)$$
$$w^{k+1} = w^k - \alpha z^{k+1}$$

# No momentum



Step-size α = 0.0038

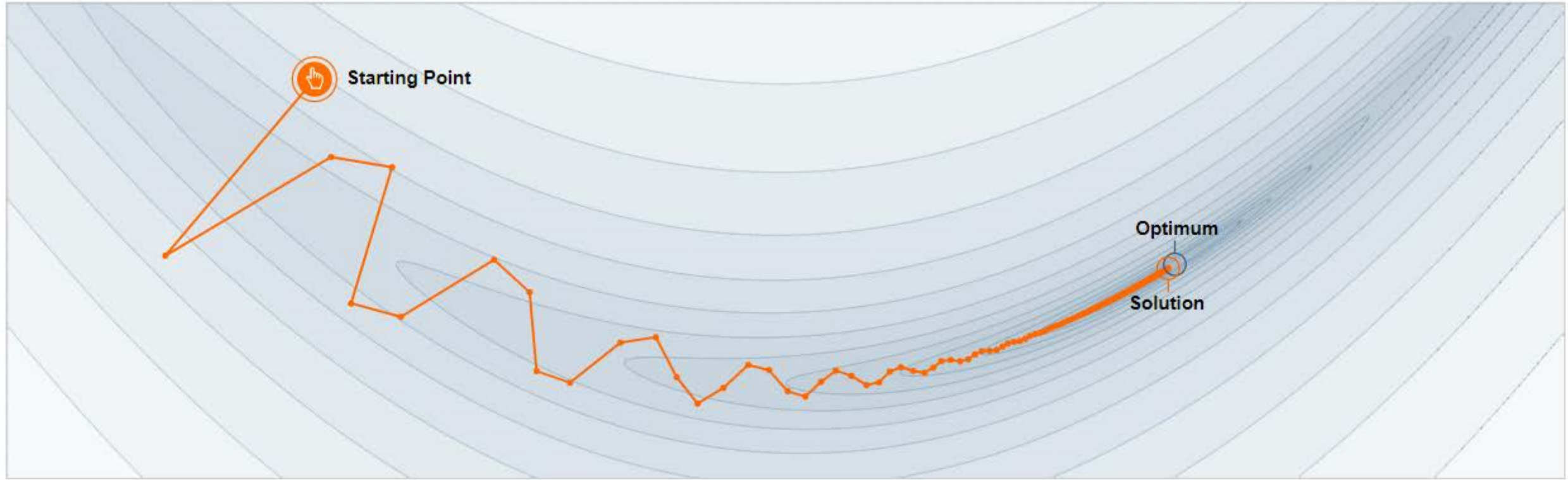Momentum β = 0.0

We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

# With appr. momentum



**Step-size α = 0.0038**

0    0.003    0.006

**Momentum β = 0.78**

0.00    0.500    0.990

We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?
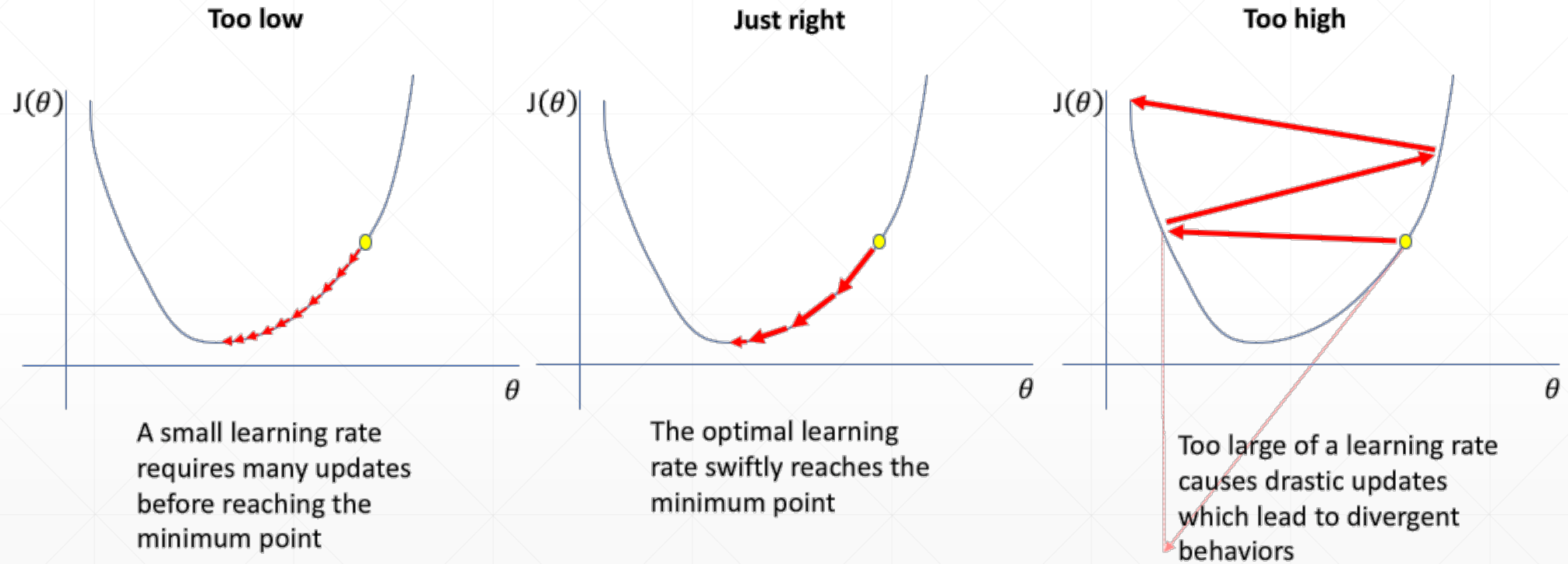
# Momentum

```
optimizer = SGD(learning_rate=0.02, momentum=0.9)
optimizer = RMSprop(learning_rate=0.02, momentum=0.9)

optimizer = SGD(learning_rate=0.02,
    beta_1=0.9,
    beta_2=0.999)
```

# Learning rate tunning

**Too low**

$J(\theta)$

$\theta$

A small learning rate requires many updates before reaching the minimum point

**Just right**

$J(\theta)$

$\theta$

The optimal learning rate swiftly reaches the minimum point

**Too high**

$J(\theta)$

$\theta$

Too large of a learning rate causes drastic updates which lead to divergent behaviors

# Learning rate decay

**Just right**

$J(\theta)$

$\theta$

$\theta$

The optimal learning rate swiftly reaches the minimum point

Step decay of learning rate

# Adaptive learning rate

```python
optimizer = SGD(learning_rate=0.2)

for epoch in range(100):
    # get loss

    # change learning rate
    optimizer.learning_rate = 0.2 * (100-epoch)/100

    # update weights
```

# 下一课时

Early Stopping, Dropout

# Thank You.