

动量与学习率

主讲：龙良曲

Outline

- momentum
 - learning rate decay
-

Momentum

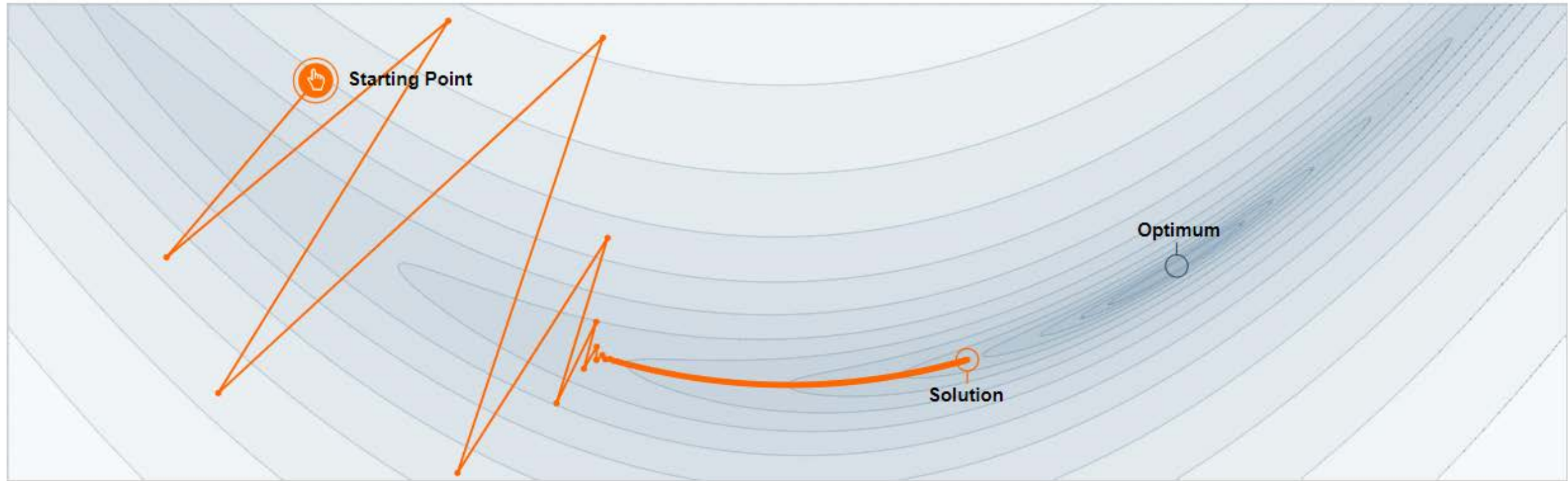
$$w^{k+1} = w^k - \alpha \nabla f(w^k).$$

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

$$w^{k+1} = w^k - \alpha z^{k+1}$$

更新方向不仅与当前梯度有关，还与前一次更新方向有关
(不容易陷入局部最优)

No momentum



Step-size $\alpha = 0.0038$

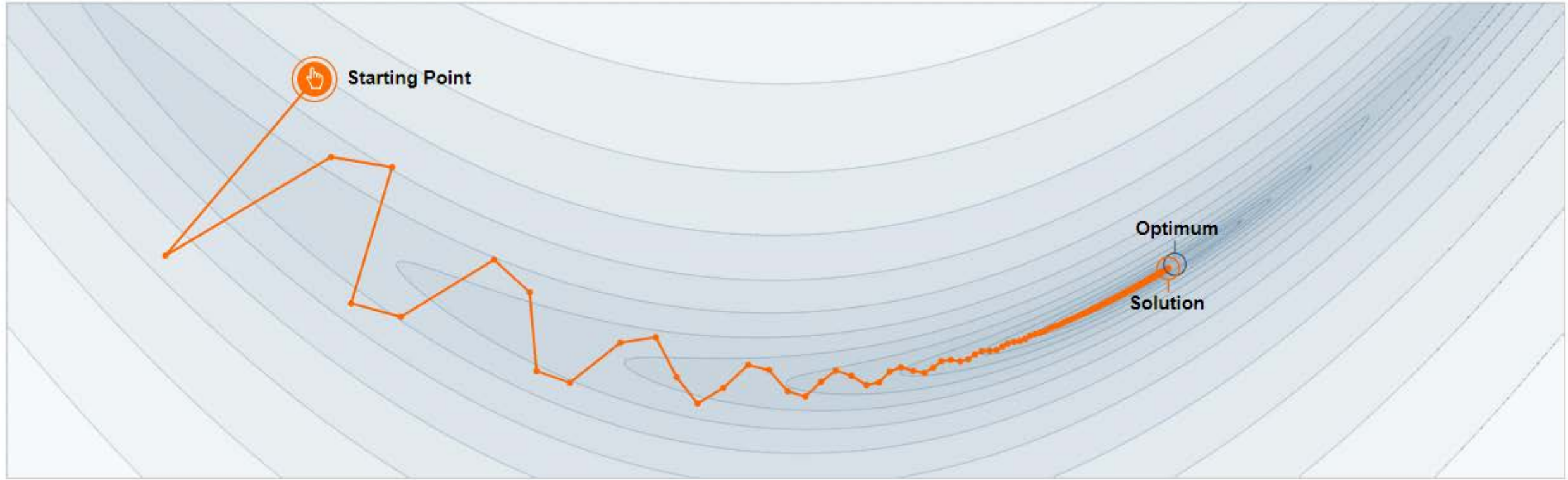


Momentum $\beta = 0.0$



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 $\alpha = 0.0038$



Momentum $\beta = 0.78$



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

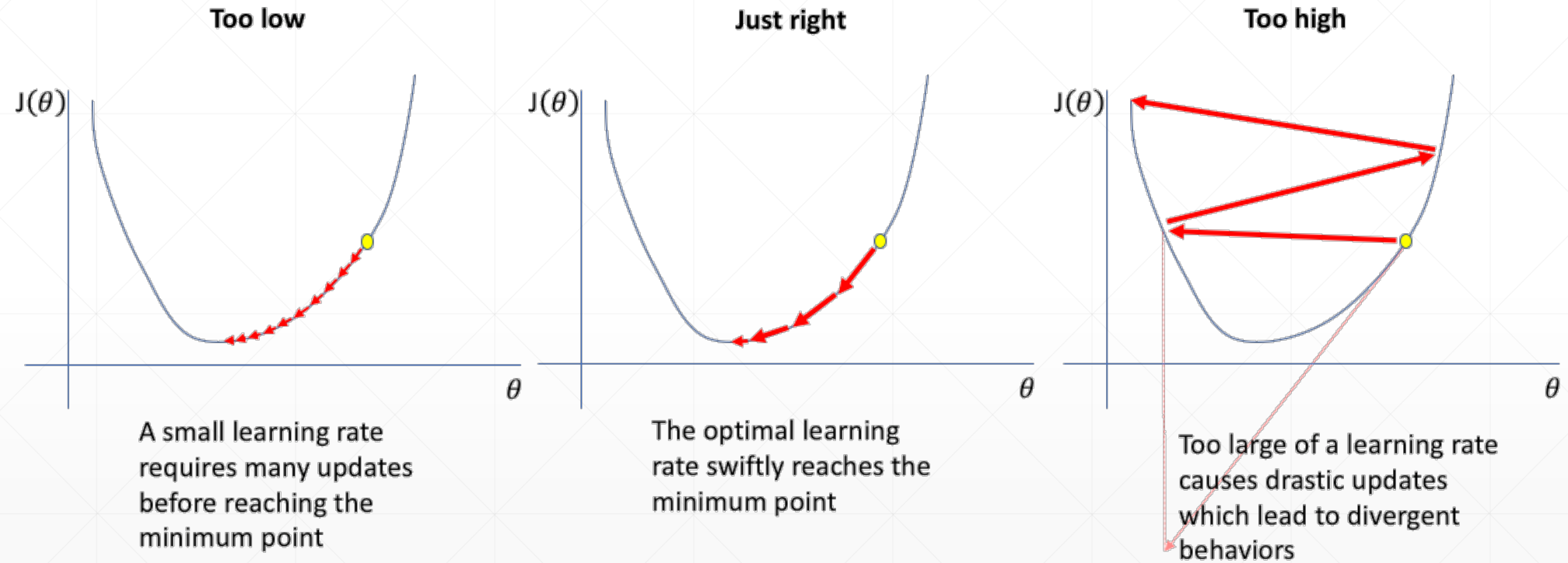


```
optimizer = SGD(learning_rate=0.02, momentum=0.9)
optimizer = RMSprop(learning_rate=0.02, momentum=0.9)

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

内含momentum优化策略

Learning rate tuning





Andrej Karpathy 

@karpathy



3e-4 is the best learning rate for Adam, hands down.

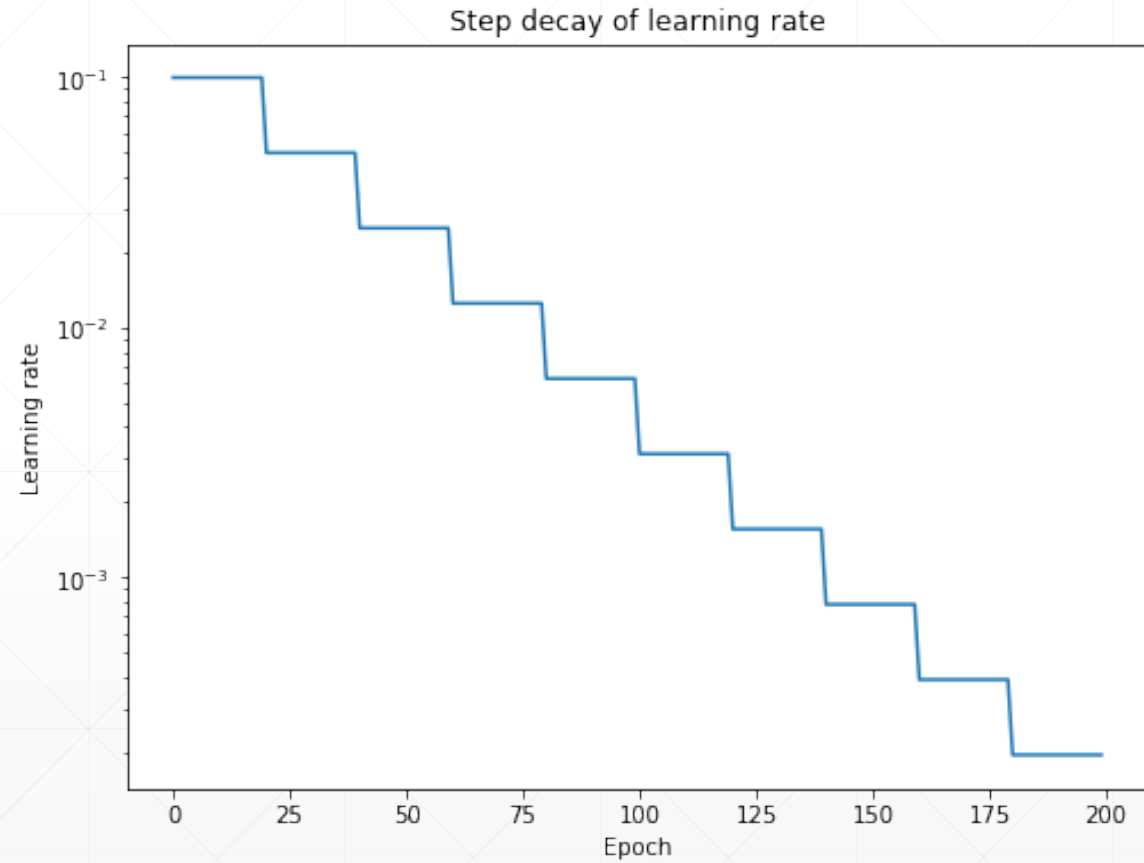
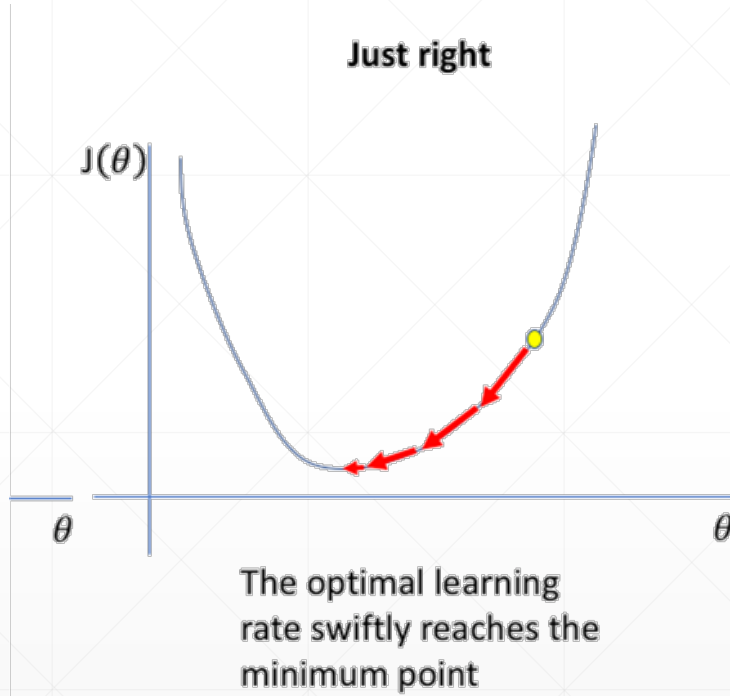
♡ 408 11:01 AM - Nov 24, 2016



💬 124 people are talking about this



Learning rate decay





Adaptive learning rate



```
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.
