



PyTorch Tutorial

03. Gradient Descent

Revision

- What would be the best model for the data?
- Linear model?

x (hours)	y (points)
1	2
2	4
3	6
4	?



Linear Model

$$\hat{y} = x * \omega$$

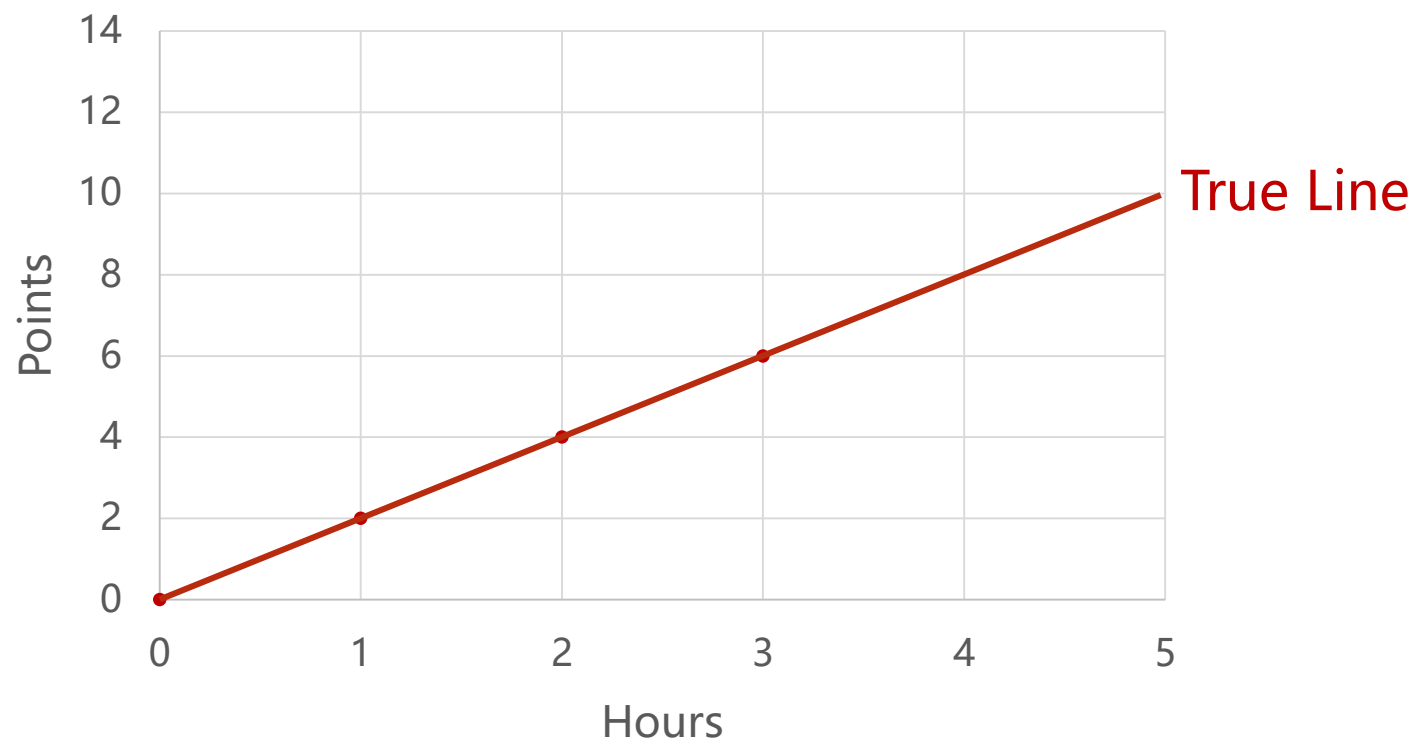
To simplify the model

Revision

Linear Model

$$\hat{y} = x * \omega$$

x (hours)	y (points)
1	2
2	4
3	6



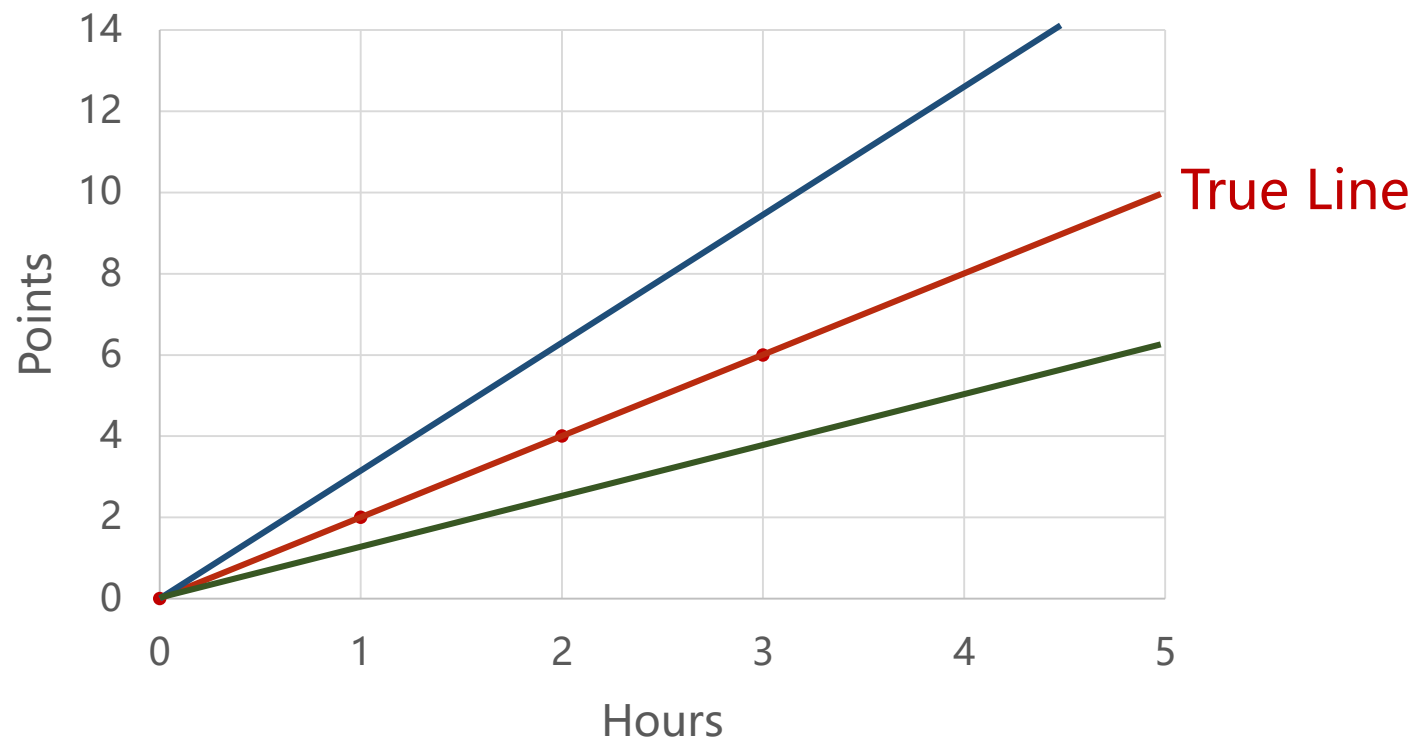
Revision

Linear Model

$$\hat{y} = x * \omega$$

x (hours)	y (points)
1	2
2	4
3	6

The machine starts with **a random guess**, ω = random value

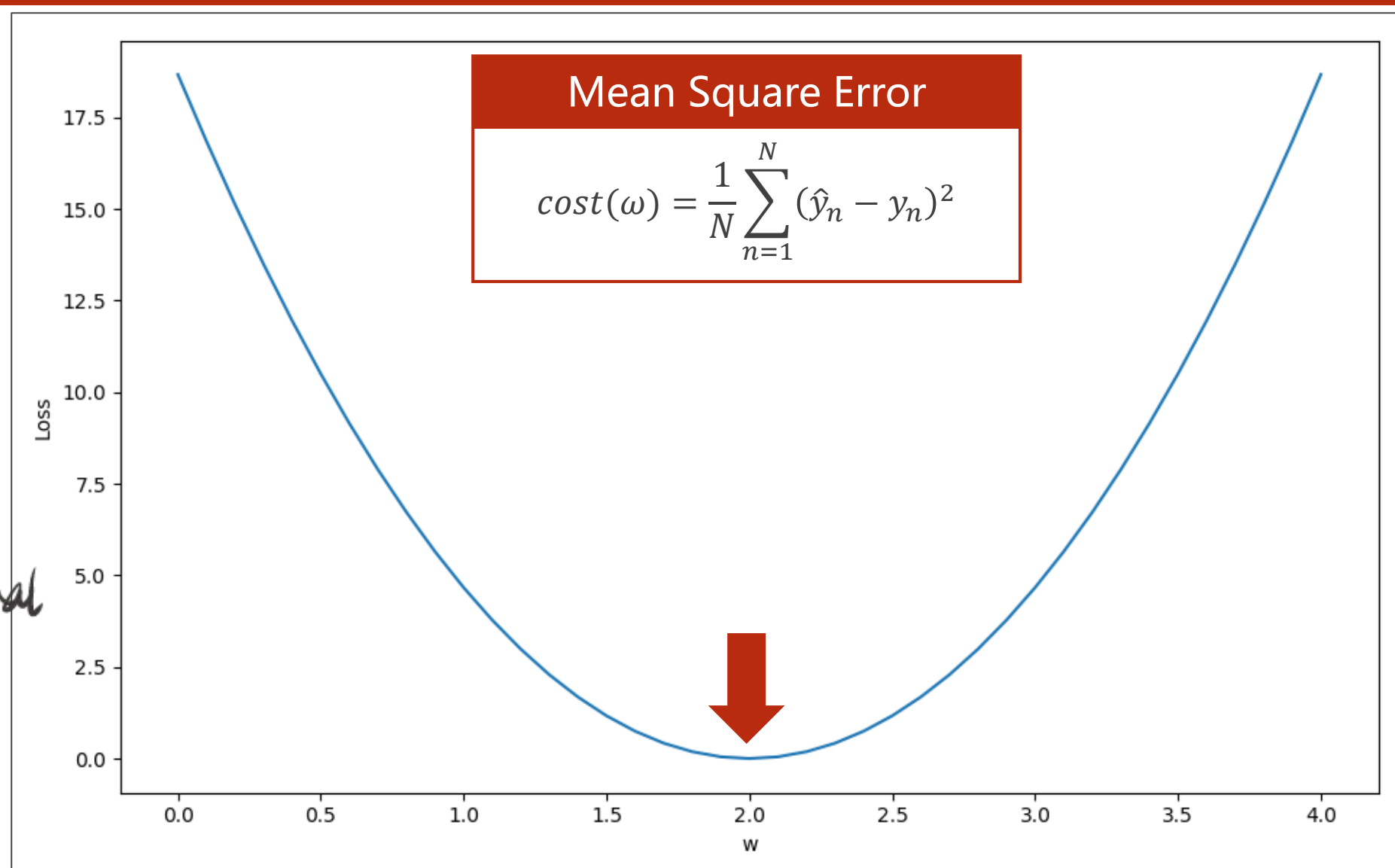


Revision

It can be found that when $\omega = 2$, the cost will be minimal.

will be minimal

肉眼直接观察到 minimal



Lecturer : Hongpu Liu

Lecture 3-5

PyTorch Tutorial @ SLAM Research Group

一、穷举法

使用穷举法列举所有可能权重 w_1
 对应的目标损失函数 $loss_i$
 看哪一种权重组合的损失值最小

仅一个权重 $w_1 \rightarrow 100$ 个数据, 搜索量为 100
 两个权重 $w_1, w_2 \rightarrow 100 \dots 100^2$
 十个权重 $w_1, w_2, \dots, w_{10} \rightarrow 100 \dots 100^{10}$

二、分治法

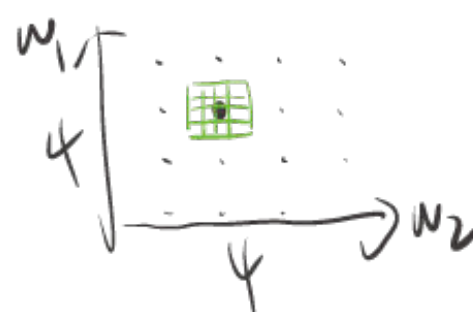
① 若搜索区间 = 100, 不直接分为 100 份 ($100 \times 100 = 10000$)

② 可先分为 4 份 $4 \times 4 = 16$, 选出 16 个点中的最小值

再以 16 个点中的最低值为区间范围, 继续划分为 4 份 $4 \times 4 = 16$

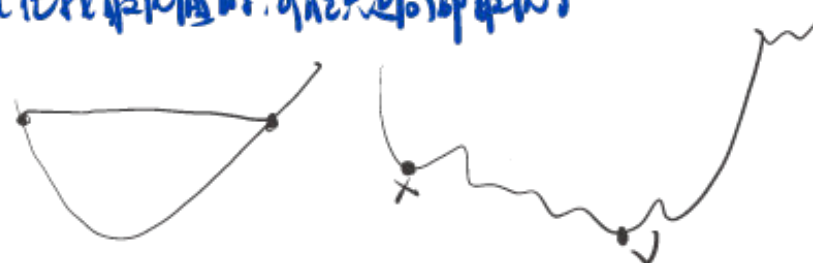
两层一共搜索 $16 + 16 = 32$

比直接遍历未访问, 否则搜索量 = $16 \times 16 = 256$

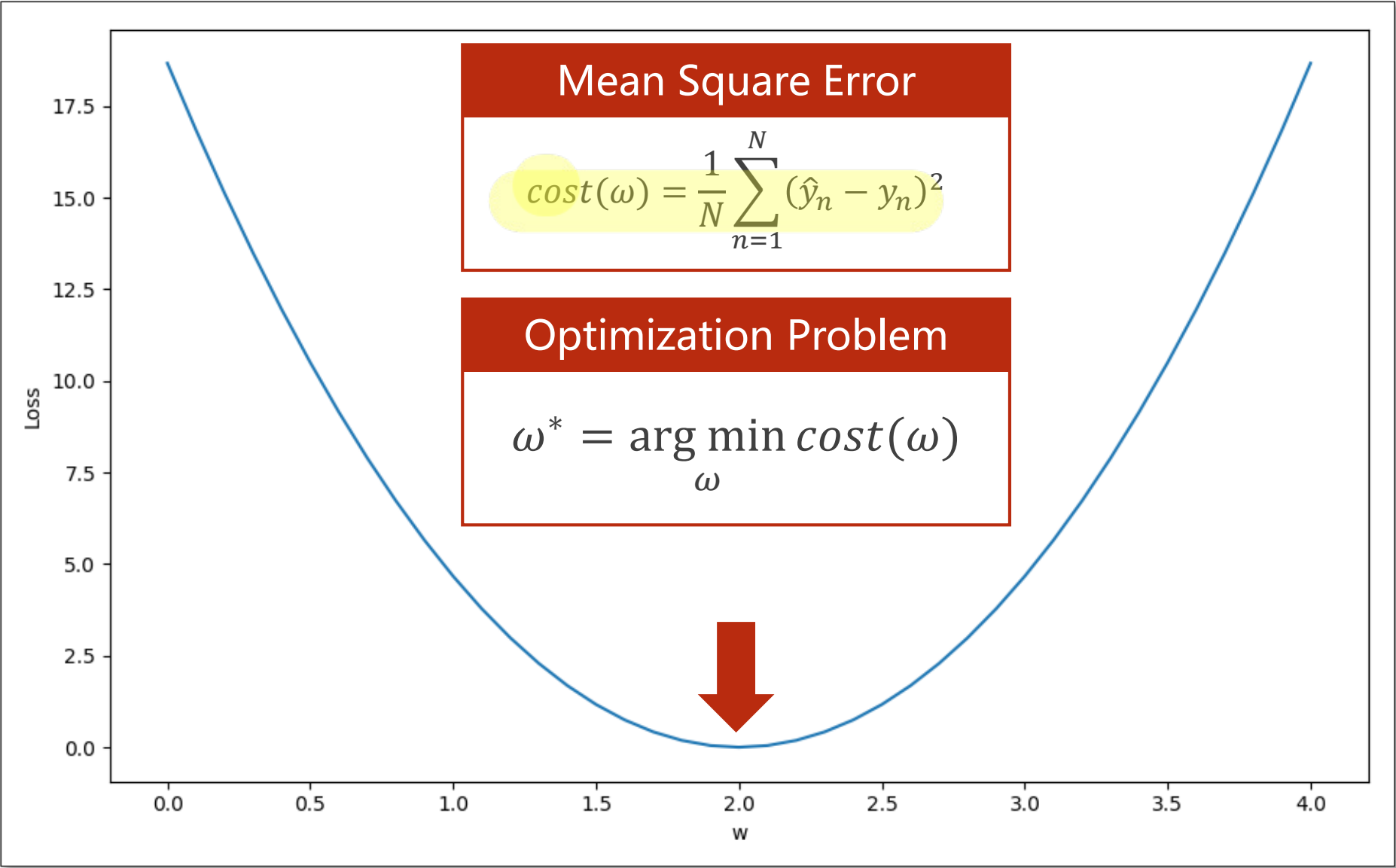


③ 分治算法可减少搜索量

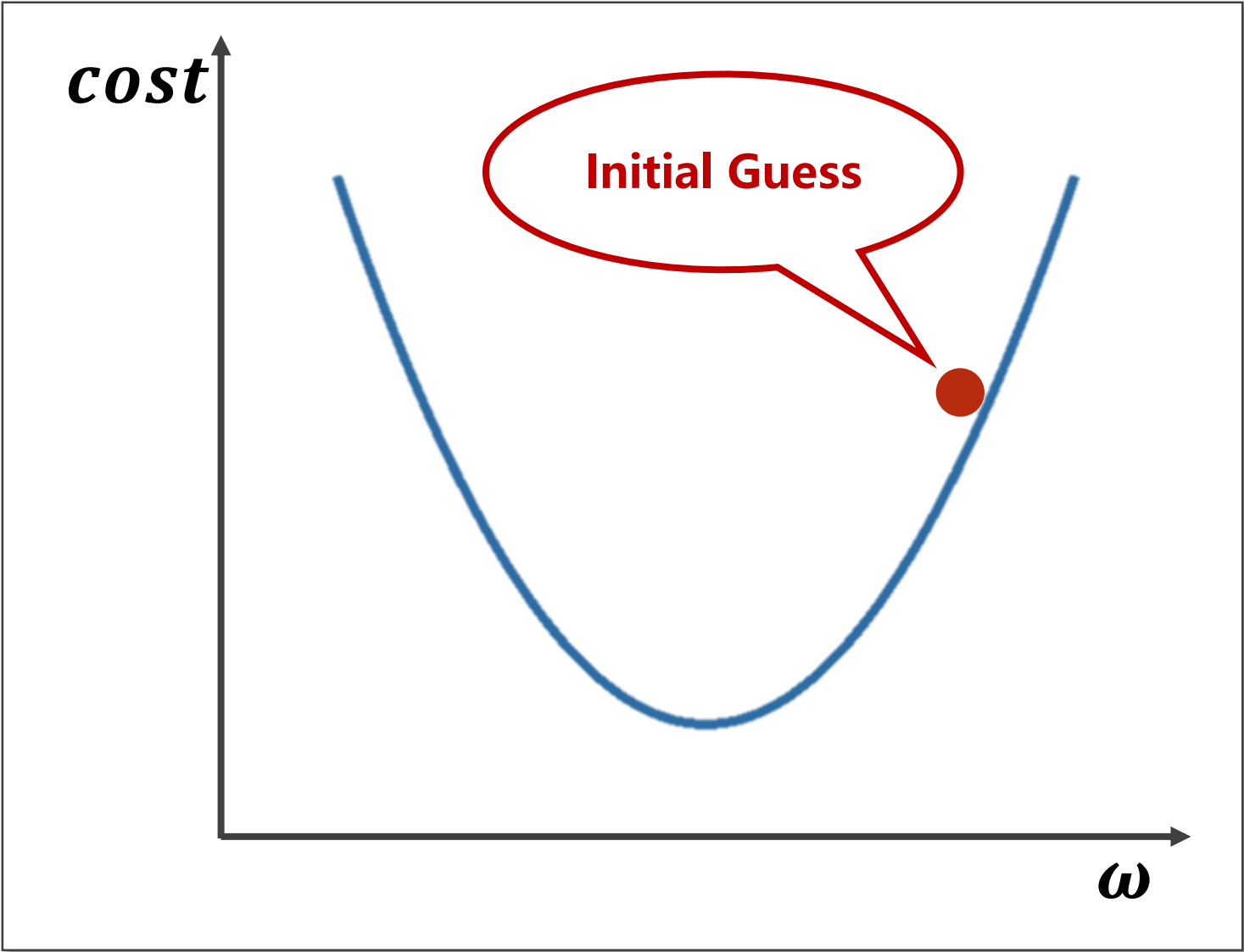
有一定优势,
 但我们的损失函数不一定是连续平滑的凸函数
 (在找最优值时, 可能只是局部最优)



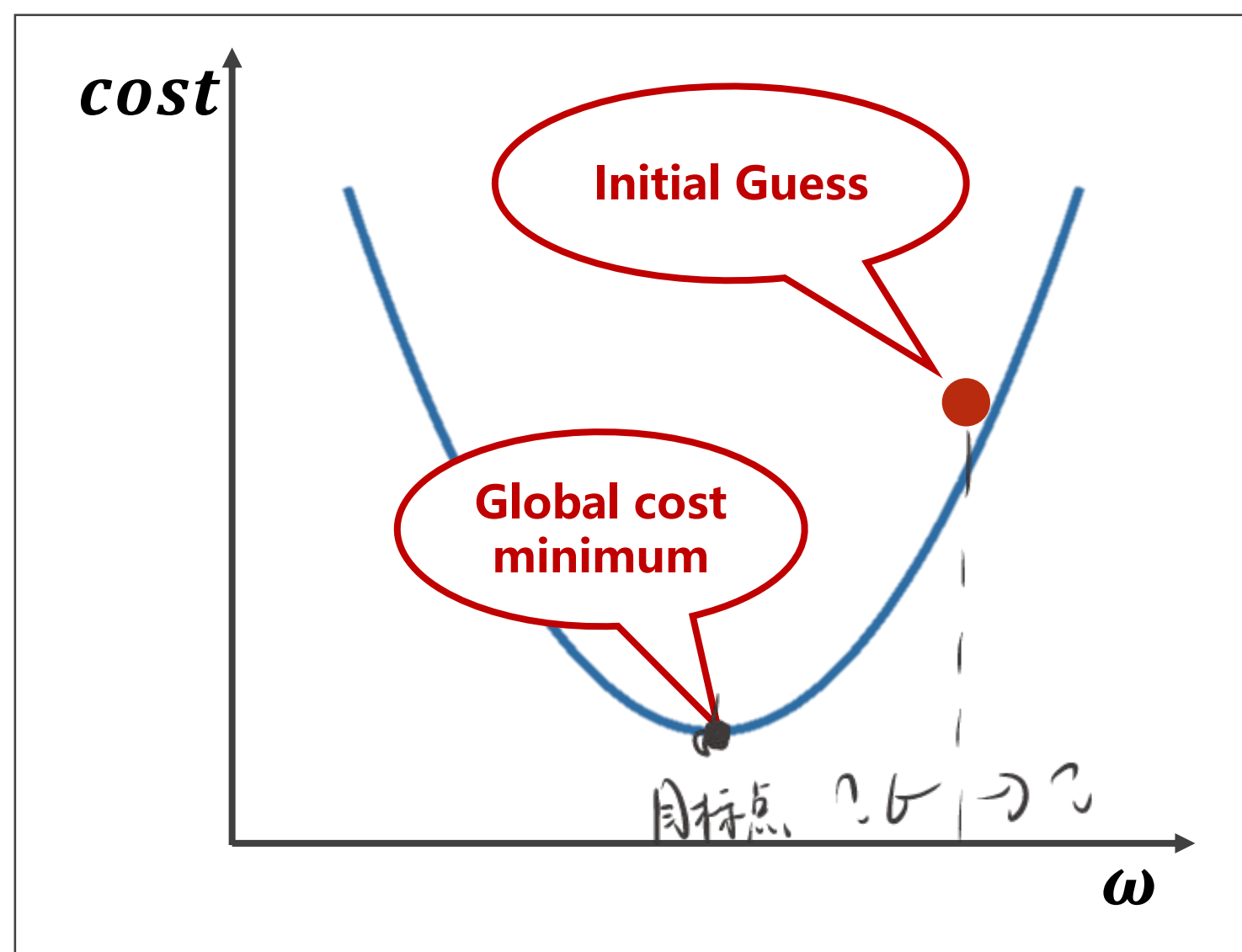
Optimization Problem



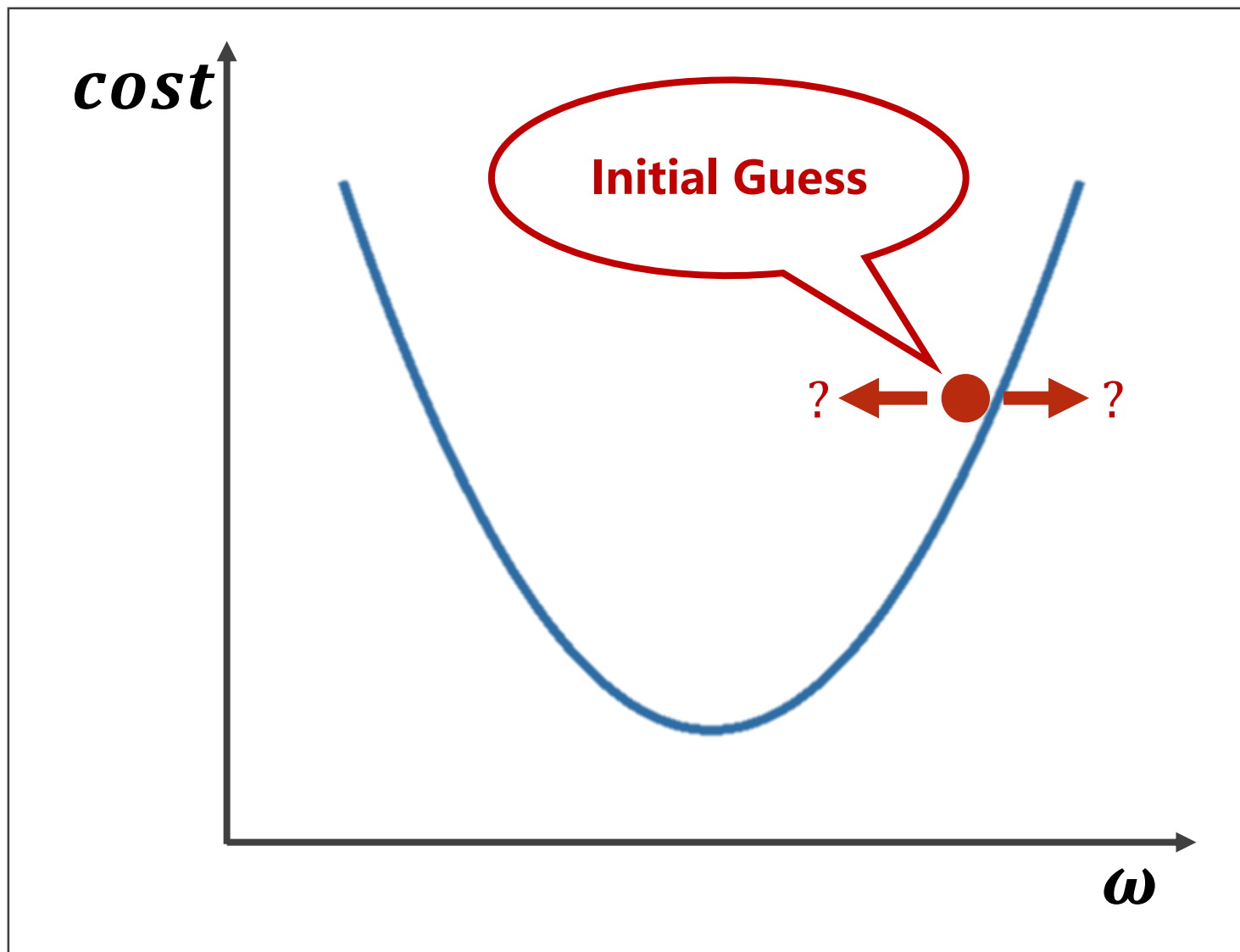
Gradient Descent Algorithm



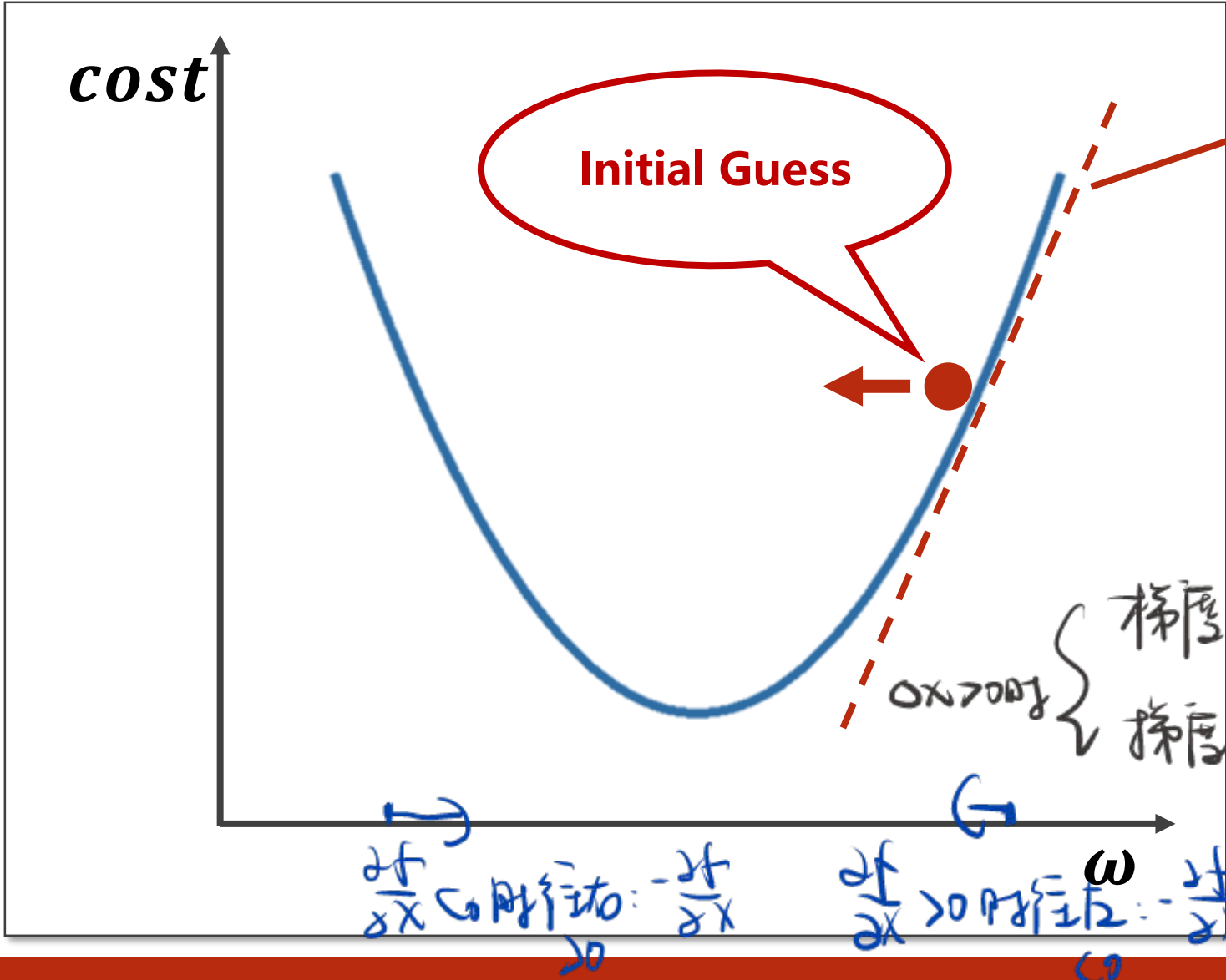
Gradient Descent Algorithm



Gradient Descent Algorithm



Gradient Descent Algorithm



Gradient

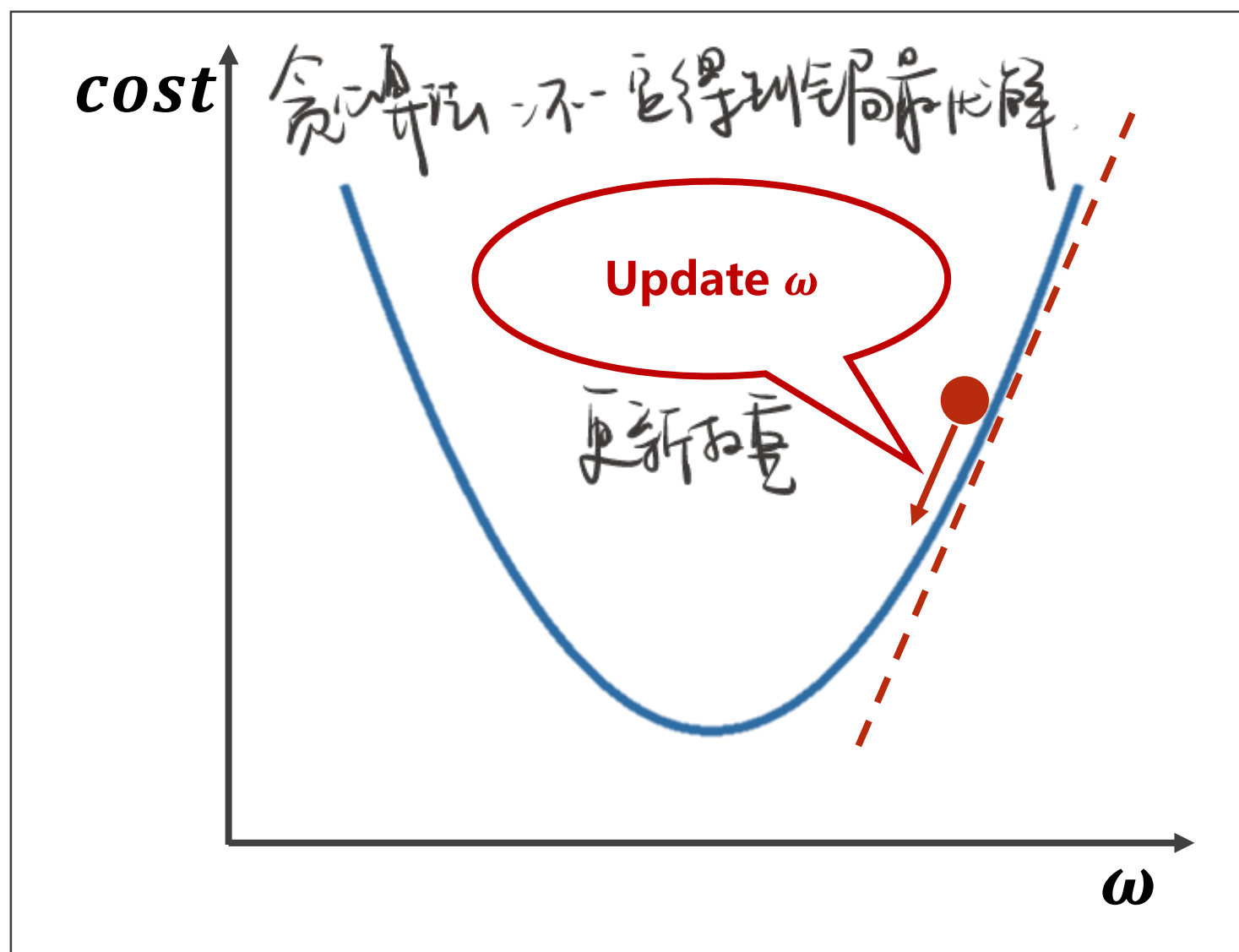
$$\frac{\partial cost}{\partial w}$$

梯度，损失函数对参数的导数 = 该点的切线斜率

$$\frac{\partial f}{\partial x} = \lim_{\Delta x \rightarrow 0} \frac{f(x+\Delta x) - f(x)}{\Delta x}$$

梯度 $\frac{\partial f}{\partial x} > 0 \Rightarrow$ 正方向 \rightarrow 上升 \Rightarrow 负方向 $-\frac{\partial f}{\partial x} \downarrow$
梯度 $\frac{\partial f}{\partial x} < 0 \Rightarrow$ 负方向 \rightarrow 下降 \Rightarrow 正方向 $-\frac{\partial f}{\partial x} \downarrow$

Gradient Descent Algorithm



Gradient

$$\frac{\partial cost}{\partial \omega}$$

Update

$$\omega = \omega - \alpha \frac{\partial cost}{\partial \omega}$$

学习率: 表示往前走多远
一般要取得小一点, 否则很难收敛

Lecturer : Hongpu Liu

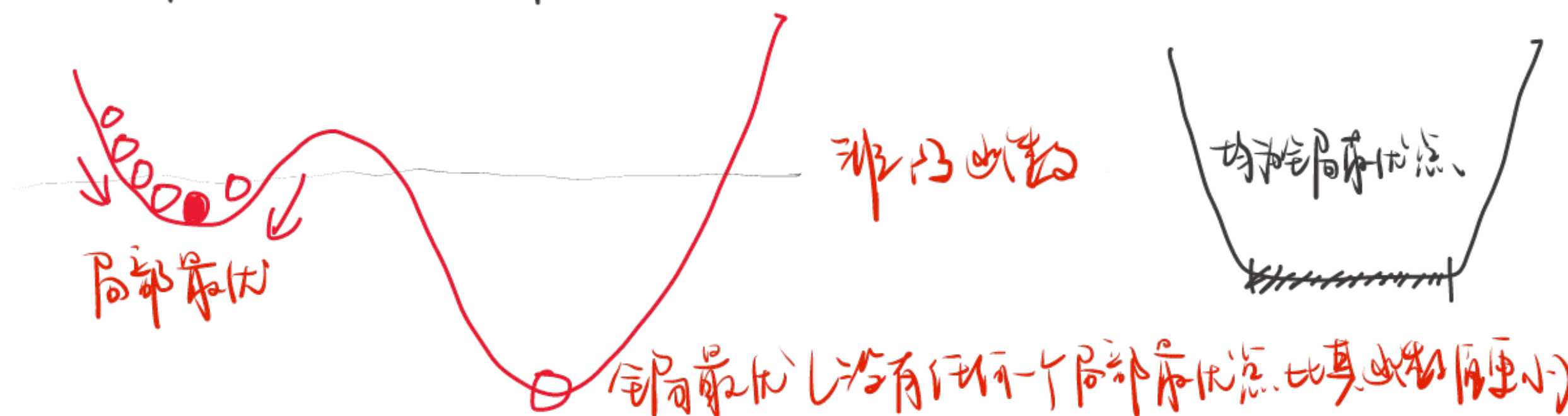
Lecture 3-11

PyTorch Tutorial @ SLAM Research Group

梯度下降: 很难找到全局最优解, 但深度学习会大量使用梯度下降作为最基本的算法

以前大家认为: 陷入局部最优解

但实际上深度神经网络中, 损失函数并没有很多局部最优解

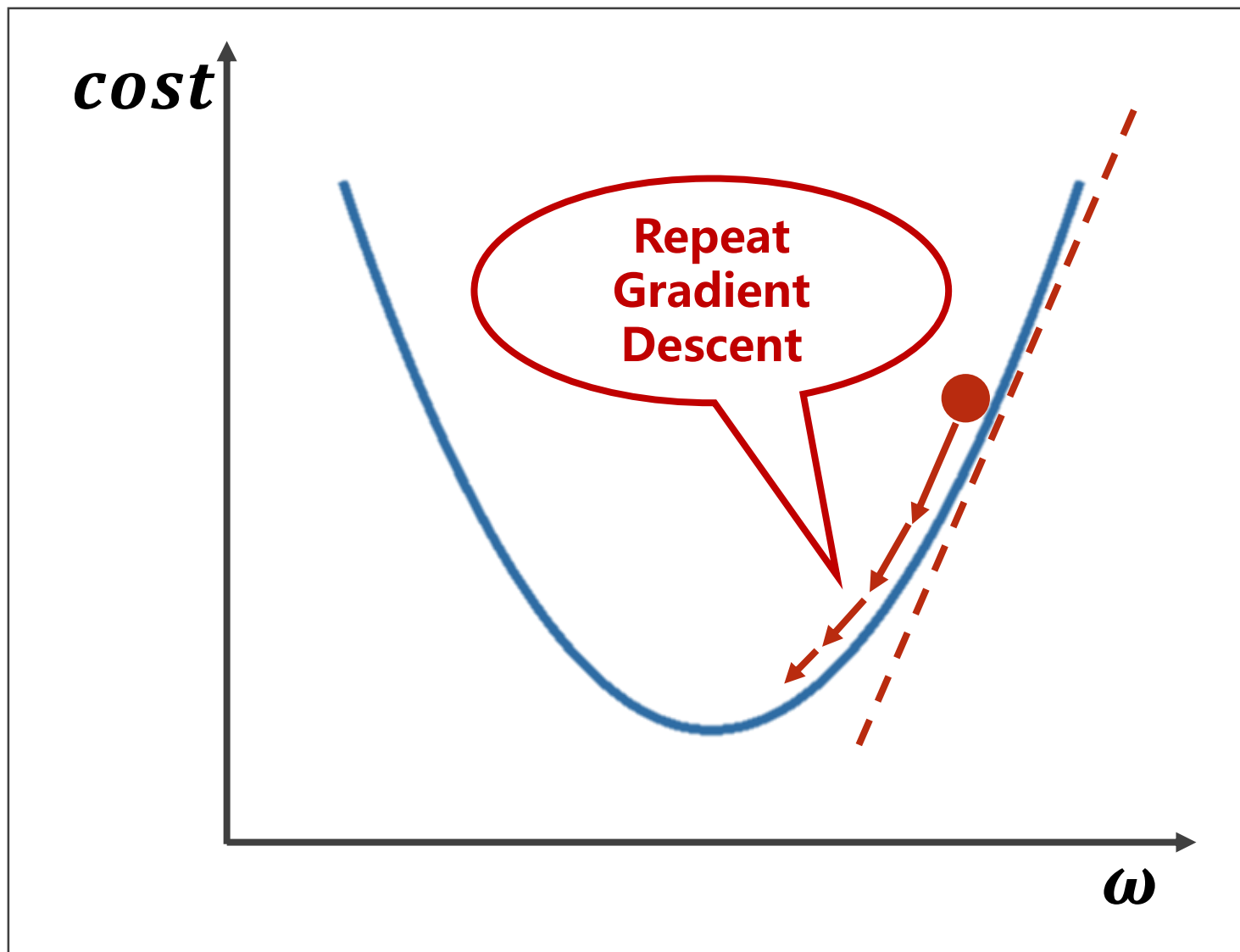


所以主要的问题是收敛速度, 什么时候迭代才停下来

$w = w - \alpha g$
 $g=0$: 每次梯度更新都无效
 ps: 该函数中没有局部最优解



Gradient Descent Algorithm



Gradient

$$\frac{\partial cost}{\partial \omega}$$

Update

$$\omega = \omega - \alpha \frac{\partial cost}{\partial \omega}$$

Gradient Descent Algorithm

Derivative

$$\frac{\partial cost(\omega)}{\partial \omega} = \frac{\partial}{\partial \omega} \frac{1}{N} \sum_{n=1}^N (x_n \cdot \omega - y_n)^2$$

和的导数=导数的和

$$= \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial \omega} (x_n \cdot \omega - y_n)^2$$

链式求导

$$= \frac{1}{N} \sum_{n=1}^N 2 \cdot (x_n \cdot \omega - y_n) \frac{\partial (x_n \cdot \omega - y_n)}{\partial \omega}$$

链式求导

$$= \frac{1}{N} \sum_{n=1}^N 2 \cdot x_n \cdot (x_n \cdot \omega - y_n)$$

$\frac{\partial z^2}{\partial z} = 2z \cdot \frac{\partial z}{\partial \omega}$

Gradient

$$\frac{\partial cost}{\partial \omega}$$

Update

$$\omega = \omega - \alpha \frac{\partial cost}{\partial \omega}$$

Gradient Descent Algorithm

Derivative

$$\begin{aligned}\frac{\partial cost(\omega)}{\partial \omega} &= \frac{\partial}{\partial \omega} \frac{1}{N} \sum_{n=1}^N (x_n \cdot \omega - y_n)^2 \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial \omega} (x_n \cdot \omega - y_n)^2 \\ &= \frac{1}{N} \sum_{n=1}^N 2 \cdot (x_n \cdot \omega - y_n) \frac{\partial (x_n \cdot \omega - y_n)}{\partial \omega} \\ &= \frac{1}{N} \sum_{n=1}^N 2 \cdot x_n \cdot (x_n \cdot \omega - y_n)\end{aligned}$$

Gradient

$$\frac{\partial cost}{\partial \omega}$$

Update

$$\omega = \omega - \alpha \frac{\partial cost}{\partial \omega}$$

Update

$$\omega = \omega - \alpha \frac{1}{N} \sum_{n=1}^N 2 \cdot x_n \cdot (x_n \cdot \omega - y_n)$$

Implementation

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]

w = 1.0

def forward(x):
    return x * w

def cost(xs, ys):
    cost = 0
    for x, y in zip(xs, ys):
        y_pred = forward(x)
        cost += (y_pred - y) ** 2
    return cost / len(xs)

def gradient(xs, ys):
    grad = 0
    for x, y in zip(xs, ys):
        grad += 2 * x * (x * w - y)
    return grad / len(xs)

print('Predict (before training)', 4, forward(4))
for epoch in range(100):
    cost_val = cost(x_data, y_data)
    grad_val = gradient(x_data, y_data)
    w -= 0.01 * grad_val
    print('Epoch:', epoch, 'w=', w, 'loss=', cost_val)
print('Predict (after training)', 4, forward(4))
```

X
Y

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]
```

Prepare the training set.

准备训练集

Implementation

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]
```

```
w = 1.0
```

```
def forward(x):
    return x * w
```

```
def cost(xs, ys):
    cost = 0
    for x, y in zip(xs, ys):
        y_pred = forward(x)
        cost += (y_pred - y) ** 2
    return cost / len(xs)
```

```
def gradient(xs, ys):
    grad = 0
    for x, y in zip(xs, ys):
        grad += 2 * x * (x * w - y)
    return grad / len(xs)
```

```
print('Predict (before training)', 4, forward(4))
for epoch in range(100):
    cost_val = cost(x_data, y_data)
    grad_val = gradient(x_data, y_data)
    w -= 0.01 * grad_val
    print('Epoch:', epoch, 'w=', w, 'loss=', cost_val)
print('Predict (after training)', 4, forward(4))
```

w = 1.0

Initial guess of weight.

初始猜测的权重

Implementation

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]
```

```
w = 1.0
```

```
def forward(x):
    return x * w
```

```
def cost(xs, ys):
    cost = 0
    for x, y in zip(xs, ys):
        y_pred = forward(x)
        cost += (y_pred - y) ** 2
    return cost / len(xs)
```

```
def gradient(xs, ys):
    grad = 0
    for x, y in zip(xs, ys):
        grad += 2 * x * (x * w - y)
    return grad / len(xs)
```

```
print('Predict (before training)', 4, forward(4))
for epoch in range(100):
    cost_val = cost(x_data, y_data)
    grad_val = gradient(x_data, y_data)
    w -= 0.01 * grad_val
    print('Epoch:', epoch, 'w=', w, 'loss=', cost_val)
print('Predict (after training)', 4, forward(4))
```

```
def forward(x):
    return x * w
```

Define the model:

Linear Model

$$\hat{y} = x * \omega$$

线性模型

Implementation

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]
```

```
w = 1.0
```

```
def forward(x):
    return x * w
```

```
def cost(xs, ys):
    cost = 0
    for x, y in zip(xs, ys):
        y_pred = forward(x)
        cost += (y_pred - y) ** 2
    return cost / len(xs)
```

```
def gradient(xs, ys):
    grad = 0
    for x, y in zip(xs, ys):
        grad += 2 * x * (x * w - y)
    return grad / len(xs)
```

```
print('Predict (before training)', 4, forward(4))
for epoch in range(100):
    cost_val = cost(x_data, y_data)
    grad_val = gradient(x_data, y_data)
    w -= 0.01 * grad_val
    print('Epoch:', epoch, 'w=', w, 'loss=', cost_val)
print('Predict (after training)', 4, forward(4))
```

```
def cost(xs, ys):
    cost = 0
    for x, y in zip(xs, ys):
        y_pred = forward(x)
        cost += (y_pred - y) ** 2
    return cost / len(xs)
```

Handwritten note: $(y - \hat{y})^2$

Define the cost function

Mean Square Error

$$cost(\omega) = \frac{1}{N} \sum_{n=1}^N (\hat{y}_n - y_n)^2$$

Handwritten note: 平均平方误差 MSE

Implementation

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]
```

```
w = 1.0
```

```
def forward(x):
    return x * w
```

```
def cost(xs, ys):
    cost = 0
    for x, y in zip(xs, ys):
        y_pred = forward(x)
        cost += (y_pred - y) ** 2
    return cost / len(xs)
```

```
def gradient(xs, ys):
    grad = 0
    for x, y in zip(xs, ys):
        grad += 2 * x * (x * w - y)
    return grad / len(xs)
```

```
print('Predict (before training)', 4, forward(4))
for epoch in range(100):
    cost_val = cost(x_data, y_data)
    grad_val = gradient(x_data, y_data)
    w -= 0.01 * grad_val
    print('Epoch:', epoch, 'w=', w, 'loss=', cost_val)
print('Predict (after training)', 4, forward(4))
```

```
def gradient(xs, ys):
    grad = 0
    for x, y in zip(xs, ys):
        grad += 2 * x * (x * w - y)
    return grad / len(xs)
```

Define the gradient function

Gradient

$$\frac{\partial cost}{\partial \omega} = \frac{1}{N} \sum_{n=1}^N \underbrace{2 \cdot x_n \cdot (x_n \cdot \omega - y_n)}_{\text{梯度}}$$

Implementation

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]

w = 1.0

def forward(x):
    return x * w

def cost(xs, ys):
    cost = 0
    for x, y in zip(xs, ys):
        y_pred = forward(x)
        cost += (y_pred - y) ** 2
    return cost / len(xs)

def gradient(xs, ys):
    grad = 0
    for x, y in zip(xs, ys):
        grad += 2 * x * (x * w - y)
    return grad / len(xs)

print('Predict (before training)', 4, forward(4))
for epoch in range(100):
    cost_val = cost(x_data, y_data)
    grad_val = gradient(x_data, y_data)
    w -= 0.01 * grad_val
    print('Epoch:', epoch, 'w=', w, 'loss=', cost_val)
print('Predict (after training)', 4, forward(4))
```

```
for epoch in range(100):
    cost_val = cost(x_data, y_data)
    grad_val = gradient(x_data, y_data)
    w -= 0.01 * grad_val
```

Do the update

Update

$$\omega = \omega - \alpha \frac{\partial \text{cost}}{\partial \omega}$$

训练过程更新权重

$\frac{\partial \text{cost}}{\partial w}$

Implementation

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]

w = 1.0

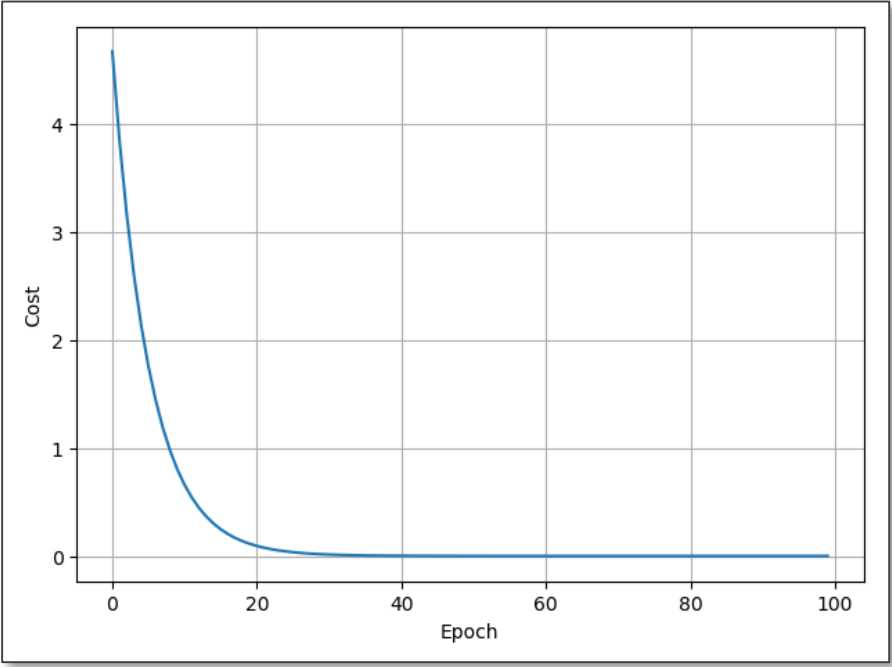
def forward(x):
    return x * w

def cost(xs, ys):
    cost = 0
    for x, y in zip(xs, ys):
        y_pred = forward(x)
        cost += (y_pred - y) ** 2
    return cost / len(xs)

def gradient(xs, ys):
    grad = 0
    for x, y in zip(xs, ys):
        grad += 2 * x * (x * w - y)
    return grad / len(xs)

print('Predict (before training)', 4, forward(4))
for epoch in range(100):
    cost_val = cost(x_data, y_data)
    grad_val = gradient(x_data, y_data)
    w -= 0.01 * grad_val
    print('Epoch:', epoch, 'w=', w, 'loss=', cost_val)
print('Predict (after training)', 4, forward(4))
```

Predict (before training) 4 4.0
Epoch: 0 w= 1.09 cost= 4.67
Epoch: 1 w= 1.18 cost= 3.84
Epoch: 2 w= 1.25 cost= 3.15
Epoch: 3 w= 1.32 cost= 2.59
Epoch: 4 w= 1.39 cost= 2.13
Epoch: 5 w= 1.44 cost= 1.75
Epoch: 6 w= 1.50 cost= 1.44
Epoch: 7 w= 1.54 cost= 1.18
Epoch: 8 w= 1.59 cost= 0.97
Epoch: 9 w= 1.62 cost= 0.80
Epoch: 10 w= 1.66 cost= 0.66
.....
Epoch: 90 w= 2.00 cost= 0.00
Epoch: 91 w= 2.00 cost= 0.00
Epoch: 92 w= 2.00 cost= 0.00
Epoch: 93 w= 2.00 cost= 0.00
Epoch: 94 w= 2.00 cost= 0.00
Epoch: 95 w= 2.00 cost= 0.00
Epoch: 96 w= 2.00 cost= 0.00
Epoch: 97 w= 2.00 cost= 0.00
Epoch: 98 w= 2.00 cost= 0.00
Epoch: 99 w= 2.00 cost= 0.00
Predict (after training) 4 8.00



Cost in each epoch

C_0, C_1, C_2, \dots
 C_0', C_1', \dots
 $C_0' = C_0$
 $C_i' = \beta C_i + (1-\beta) C_{i-1}$

指数加权均值：使训练曲线更平滑
训练误差：表示训练误差
可能梯度太大
但不意味着最低点就是损失最小值。

Stochastic Gradient Descent

Gradient Descent

$$\omega = \omega - \alpha \frac{\partial cost}{\partial \omega}$$

Derivative of Cost Function

$$\frac{\partial cost}{\partial \omega} = \frac{1}{N} \sum_{n=1}^N 2 \cdot x_n \cdot (x_n \cdot \omega - y_n)$$

所有损失的平均值 cost \rightarrow 随机一个样本，损失值 loss (每个样本权重求导并更新)

$\vec{0}$ \rightarrow 因为数据是平衡的，引入一个随机噪声，使处于稳定点的情况，向两侧推动

Stochastic Gradient Descent

$$\omega = \omega - \alpha \frac{\partial loss}{\partial \omega}$$

Derivative of Loss Function

$$\frac{\partial loss_n}{\partial \omega} = 2 \cdot x_n \cdot (x_n \cdot \omega - y_n)$$

Implementation of SGD

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]

w = 1.0

def forward(x):
    return x * w

def loss(x, y):
    y_pred = forward(x)
    return (y_pred - y) ** 2

def gradient(x, y):
    return 2 * x * (x * w - y)

print('Predict (before training)', 4, forward(4))

for epoch in range(100):
    for x, y in zip(x_data, y_data):
        grad = gradient(x, y)
        w = w - 0.01 * grad
        print("\tgrad: ", x, y, grad)
        l = loss(x, y)

    print("progress:", epoch, "w=", w, "loss=", l)

print('Predict (after training)', 4, forward(4))
```

cost \rightarrow loss
求和 \rightarrow 随机

```
def loss(x, y):
    y_pred = forward(x)
    return (y_pred - y) ** 2
```

Calculate loss function:

Loss Function

$$loss = (\hat{y} - y)^2 = (x * \omega - y)^2$$

Implementation of SGD

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]

w = 1.0

def forward(x):
    return x * w

def loss(x, y):
    y_pred = forward(x)
    return (y_pred - y) ** 2

def gradient(x, y):
    return 2 * x * (x * w - y)

print('Predict (before training)', 4, forward(4))

for epoch in range(100):
    for x, y in zip(x_data, y_data):
        grad = gradient(x, y)
        w = w - 0.01 * grad
        print("\tgrad: ", x, y, grad)
        l = loss(x, y)

    print("progress:", epoch, "w=", w, "loss=", l)

print('Predict (after training)', 4, forward(4))
```

梯度
求和 → 平均

```
def gradient(x, y):
    return 2 * x * (x * w - y)
```

Calculate loss function:

Derivative of Loss Function

$$\frac{\partial loss_n}{\partial \omega} = 2 \cdot x_n \cdot (x_n \cdot \omega - y_n)$$

Implementation of SGD

```
x_data = [1.0, 2.0, 3.0]
y_data = [2.0, 4.0, 6.0]

w = 1.0

def forward(x):
    return x * w

def loss(x, y):
    y_pred = forward(x)
    return (y_pred - y) ** 2

def gradient(x, y):
    return 2 * x * (x * w - y)

print('Predict (before training)', 4, forward(4))
```

```
for epoch in range(100):
    for x, y in zip(x_data, y_data):
        grad = gradient(x, y)
        w = w - 0.01 * grad
        print("\tgrad: ", x, y, grad)
        l = loss(x, y)
```

```
    print("progress:", epoch, "w=", w, "loss=", l)

print('Predict (after training)', 4, forward(4))
```

```
for epoch in range(100):
    for x, y in zip(x_data, y_data):
        grad = gradient(x, y)
        w = w - 0.01 * grad
        print("\tgrad: ", x, y, grad)
        l = loss(x, y)
```

Update weight by every grad of sample of train set.

→ 对每个样本都分别计算梯度并更新

梯度下降 = 效率更高 (并行计算 $f(x_i) - f(x_{i+1})$, 所有无相互依赖关系)

随机梯度下降: 性能更高, 更能找到最优值 (两个样本的梯度下降是相互依赖的, 所以不可并行)

中和

$$w = \alpha \cdot \frac{\partial g(x_i)}{\partial w}$$

$$\downarrow$$

$$w = \alpha \cdot \frac{\partial g(x_{i+1})}{\partial w}$$

梯度下降: 批量的梯度下降

Batch

性能 低

时间复杂度 低 ✓

mini-batch

高 ✓

高

PyTorch Tutorial

03. Gradient Descent