

Gradient Descent

顏安孜

azyen@nycu.edu.tw

The following slides are selected from the course material of Machine Learning by Prof.
Hung-Yi Lee

Review: Gradient Descent

- In step 3, we have to solve the following optimization problem:

$$\theta^* = \arg \min_{\theta} L(\theta) \quad L: \text{loss function} \quad \theta: \text{parameters}$$

Suppose that θ has two variables $\{\theta_1, \theta_2\}$

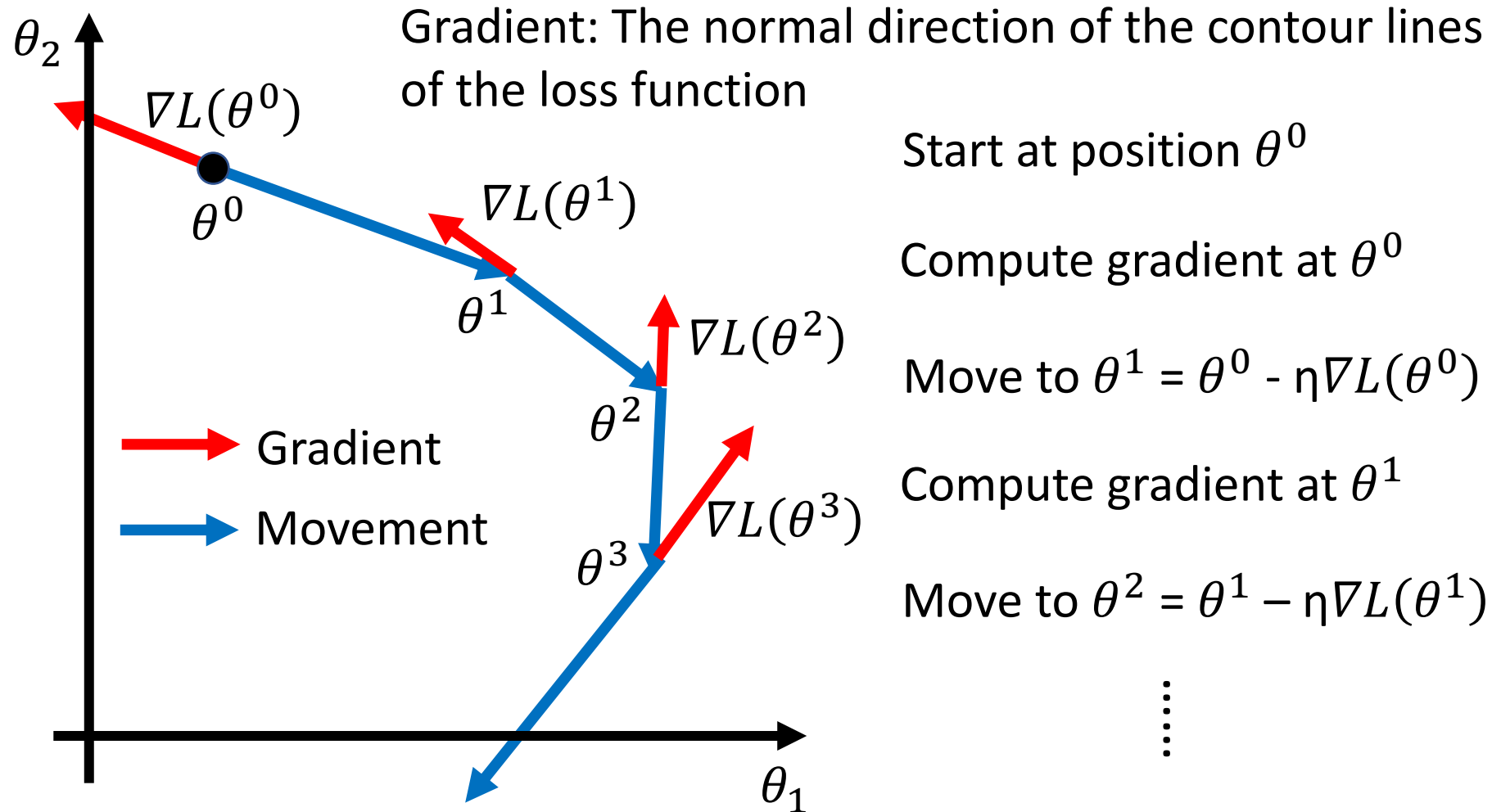
Randomly start at $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$

$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta_1)/\partial \theta_1 \\ \partial L(\theta_2)/\partial \theta_2 \end{bmatrix}$$

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta_1^0)/\partial \theta_1 \\ \partial L(\theta_2^0)/\partial \theta_2 \end{bmatrix} \Rightarrow \theta^1 = \theta^0 - \eta \nabla L(\theta^0)$$

$$\begin{bmatrix} \theta_1^2 \\ \theta_2^2 \end{bmatrix} = \begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta_1^1)/\partial \theta_1 \\ \partial L(\theta_2^1)/\partial \theta_2 \end{bmatrix} \Rightarrow \theta^2 = \theta^1 - \eta \nabla L(\theta^1)$$

Review: Gradient Descent

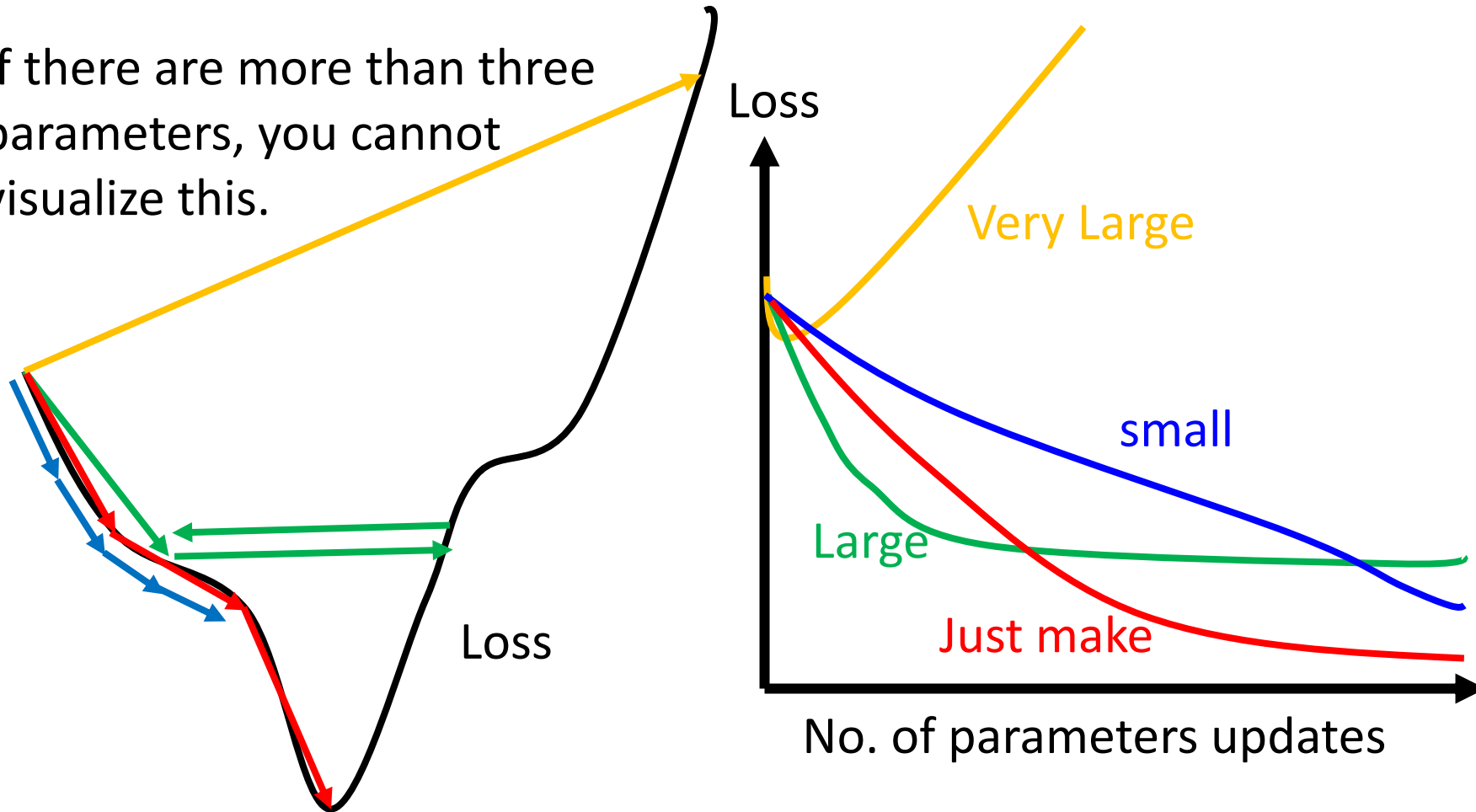


Learning Rate

$$\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$$

Set the learning rate η carefully

If there are more than three parameters, you cannot visualize this.



But you can always visualize this.

Adaptive Learning Rates

- Popular & Simple Idea: Reduce the learning rate by some factor every few epochs.
 - At the beginning, we are far from the destination, so we use larger learning rate
 - After several epochs, we are close to the destination, so we reduce the learning rate
 - E.g. 1/t decay: $\eta^t = \eta / \sqrt{t + 1}$
- Learning rate cannot be one-size-fits-all
 - Giving different parameters different learning rates

Adagrad

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad g^t = \frac{\partial L(\theta^t)}{\partial w}$$

- Divide the learning rate of each parameter by the ***root mean square of its previous derivatives***

Vanilla Gradient descent

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

w is one parameters

Adagrad

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

σ^t : ***root mean square*** of the previous derivatives of parameter w

Parameter dependent

Adagrad

σ^t : *root mean square* of the previous derivatives of parameter w

$$w^1 \leftarrow w^0 - \frac{\eta^0}{\sigma^0} g^0$$

$$w^2 \leftarrow w^1 - \frac{\eta^1}{\sigma^1} g^1$$

$$w^3 \leftarrow w^2 - \frac{\eta^2}{\sigma^2} g^2$$

\vdots

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

$$\sigma^0 = \sqrt{(g^0)^2}$$

$$\sigma^1 = \sqrt{\frac{1}{2} [(g^0)^2 + (g^1)^2]}$$

$$\sigma^2 = \sqrt{\frac{1}{3} [(g^0)^2 + (g^1)^2 + (g^2)^2]}$$

$$\sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$$

Adagrad

- Divide the learning rate of each parameter by the **root mean square of its previous derivatives**

The diagram illustrates the Adagrad update rule. It shows the parameter update equation with annotations for the learning rate and the root mean square of previous derivatives.

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

Annotations:

- $\eta^t = \frac{\eta}{\sqrt{t+1}}$ (1/t decay)
- $\sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$

A large blue arrow points from the first equation to the simplified version below:

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

Contradiction?

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad g^t = \frac{\partial L(\theta^t)}{\partial w}$$

Vanilla Gradient descent

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

Larger gradient,
larger step

Adagrad

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} \underline{g^t}$$

Larger gradient,
larger step

Larger gradient,
smaller step

Intuitive Reason

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad g^t = \frac{\partial L(\theta^t)}{\partial w}$$

- How surprise it is

g^0	g^1	g^2	g^3	g^4
0.001	0.001	0.003	0.002	0.1
g^0	g^1	g^2	g^3	g^4
10.8	20.9	31.7	12.1	0.1

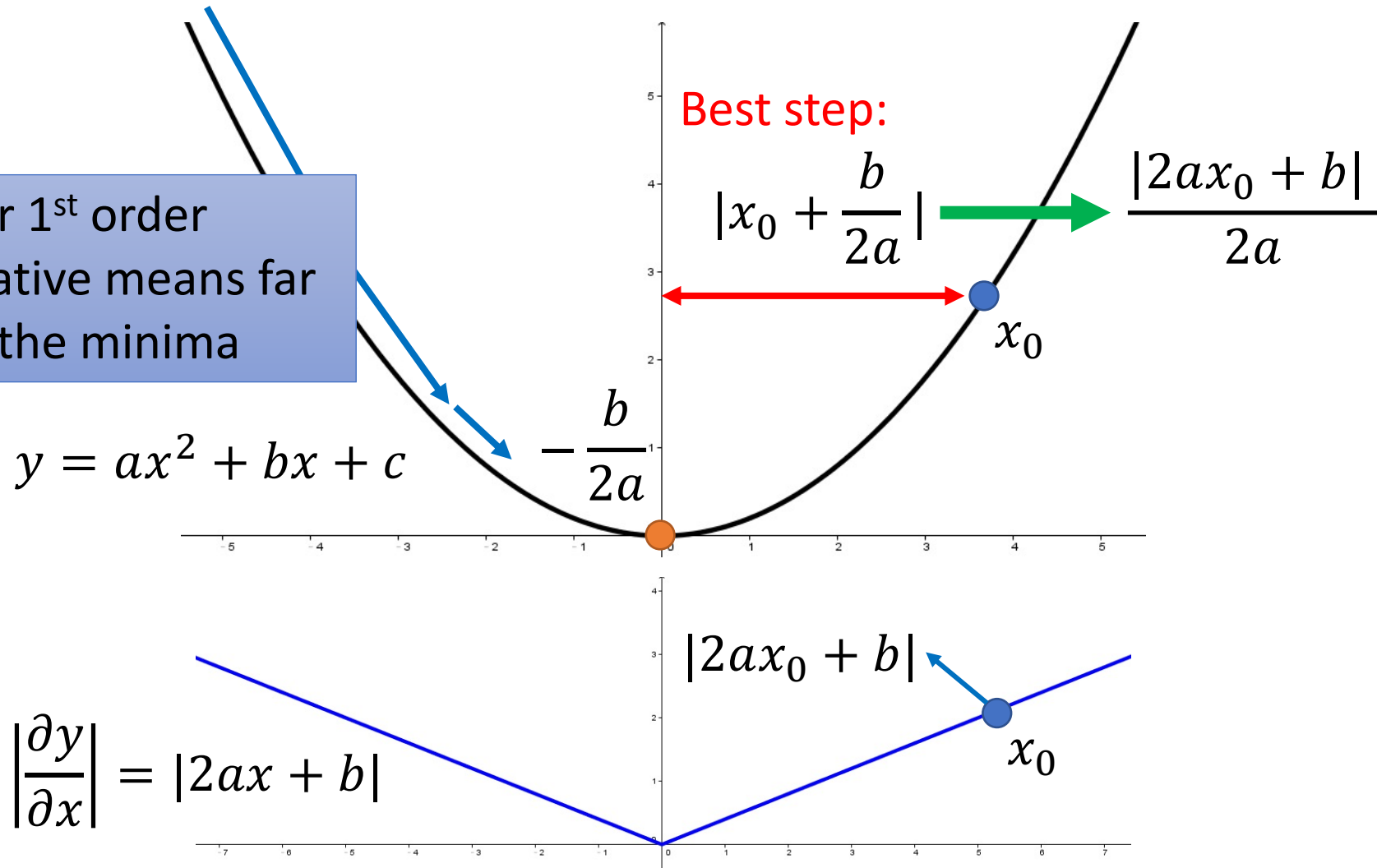
Very big

Very small

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

Larger gradient, larger steps?

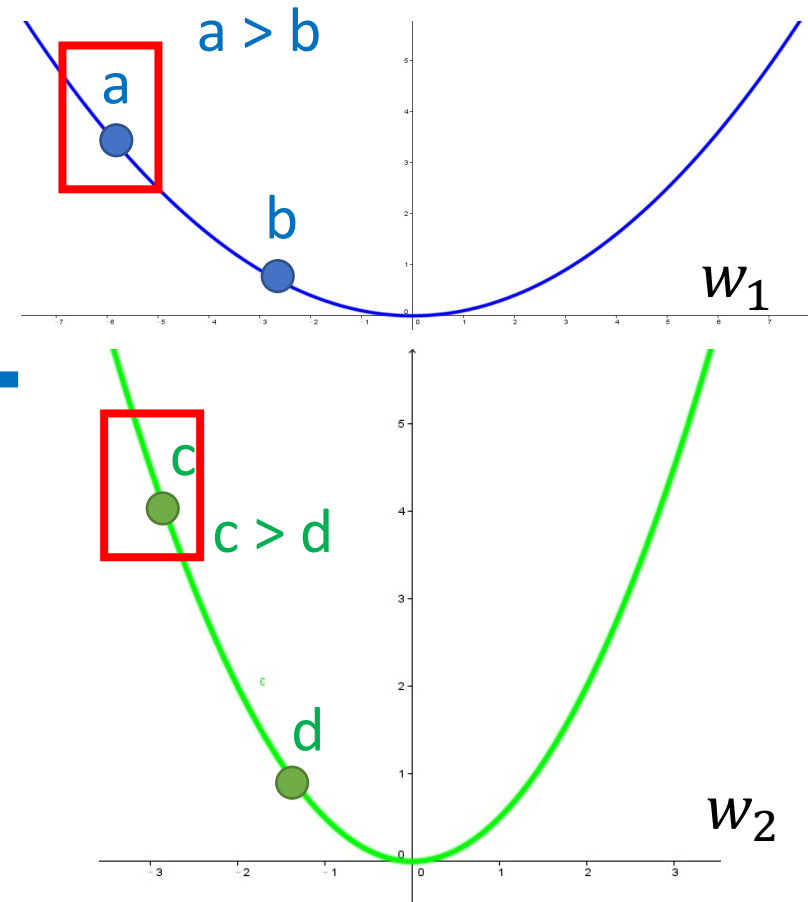
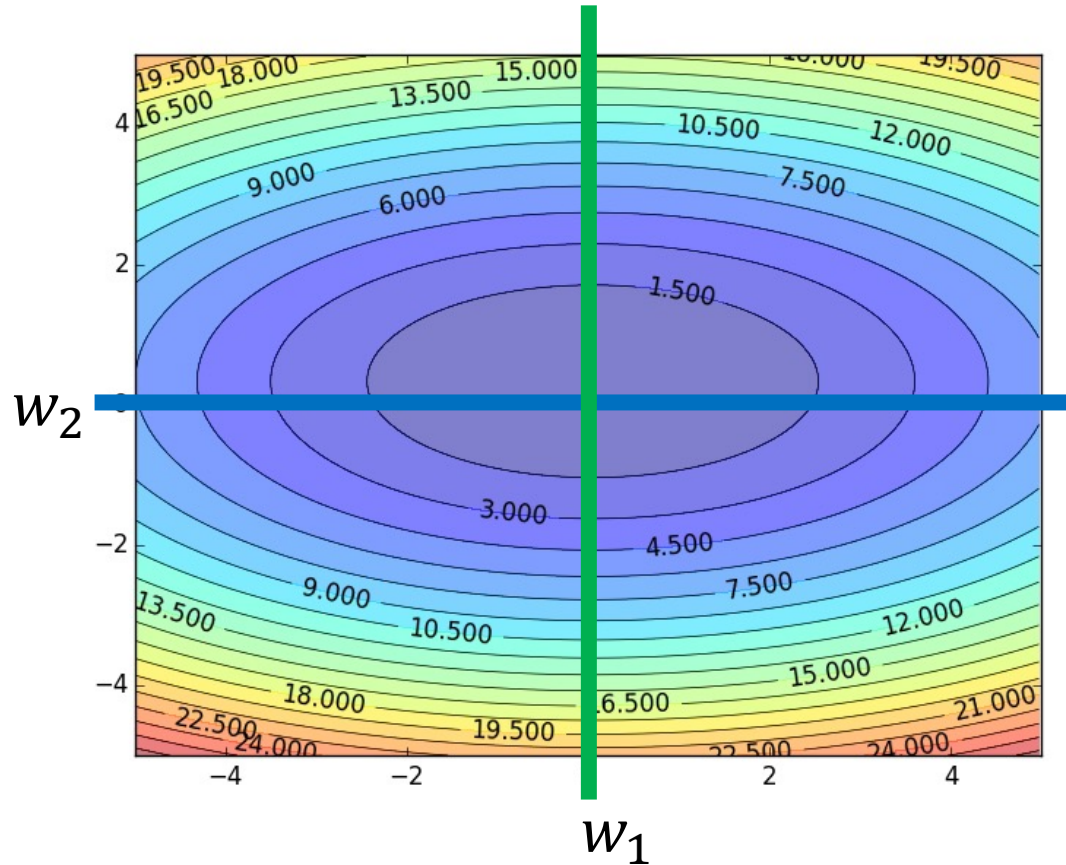
Larger 1st order derivative means far from the minima



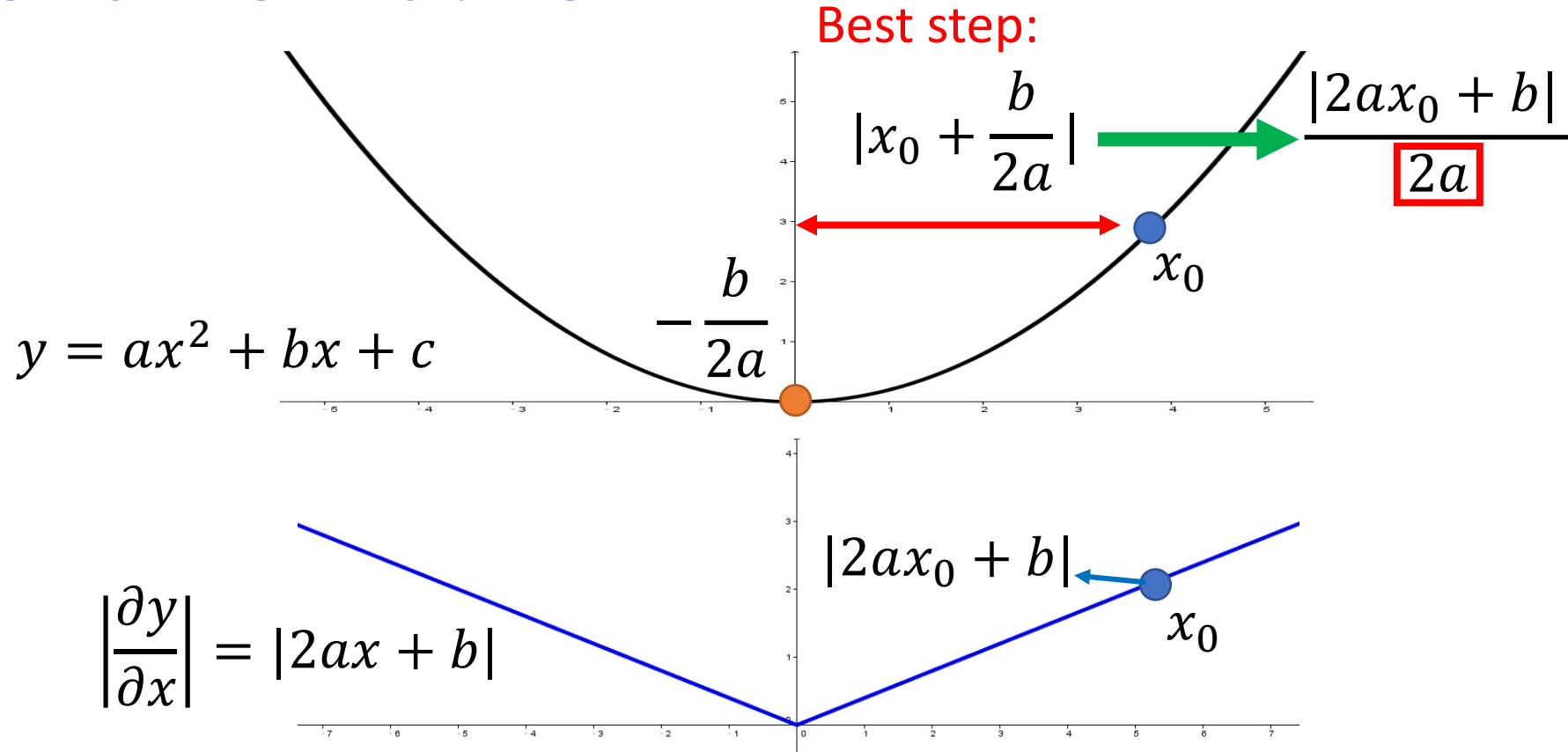
Comparison between different parameters

~~Larger 1st order derivative means far from the minima~~

Do not cross parameters



Second Derivative

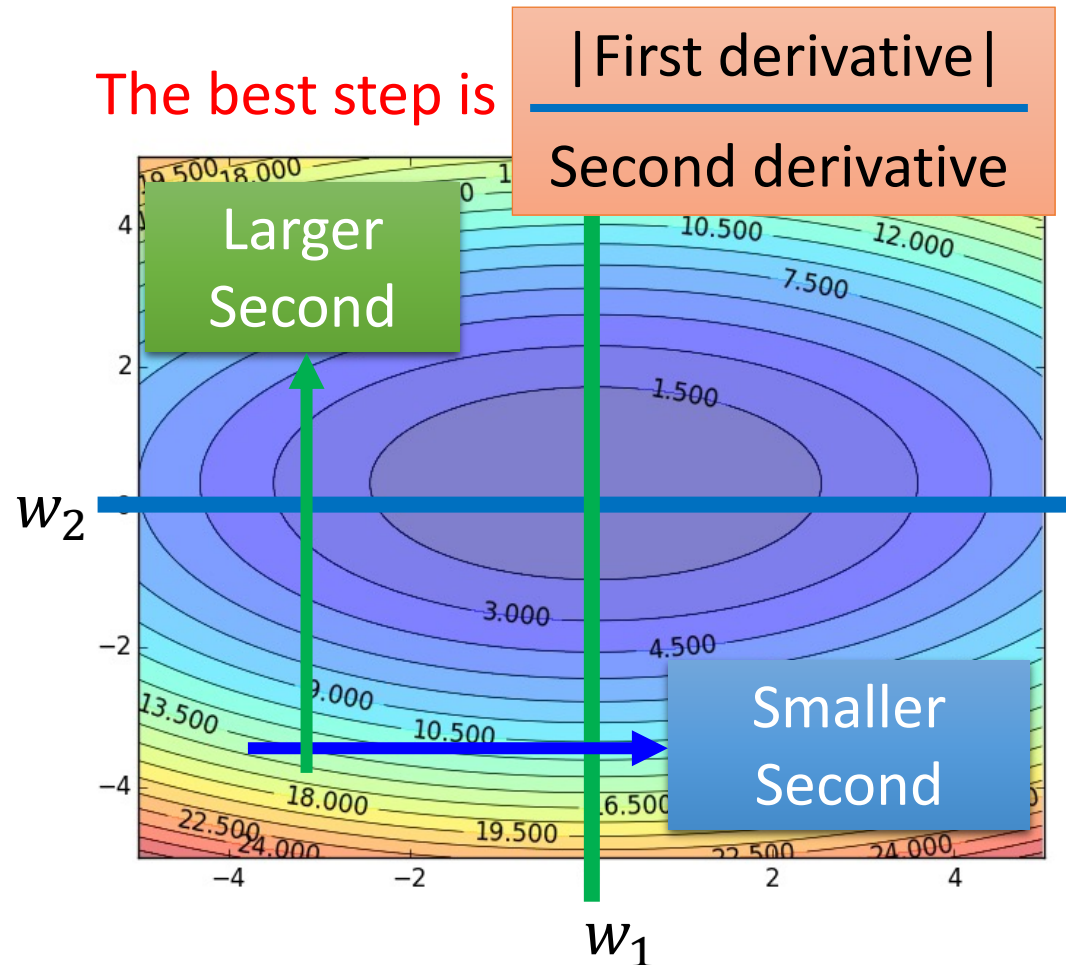


$$\frac{\partial^2 y}{\partial x^2} = 2a$$

The best step is

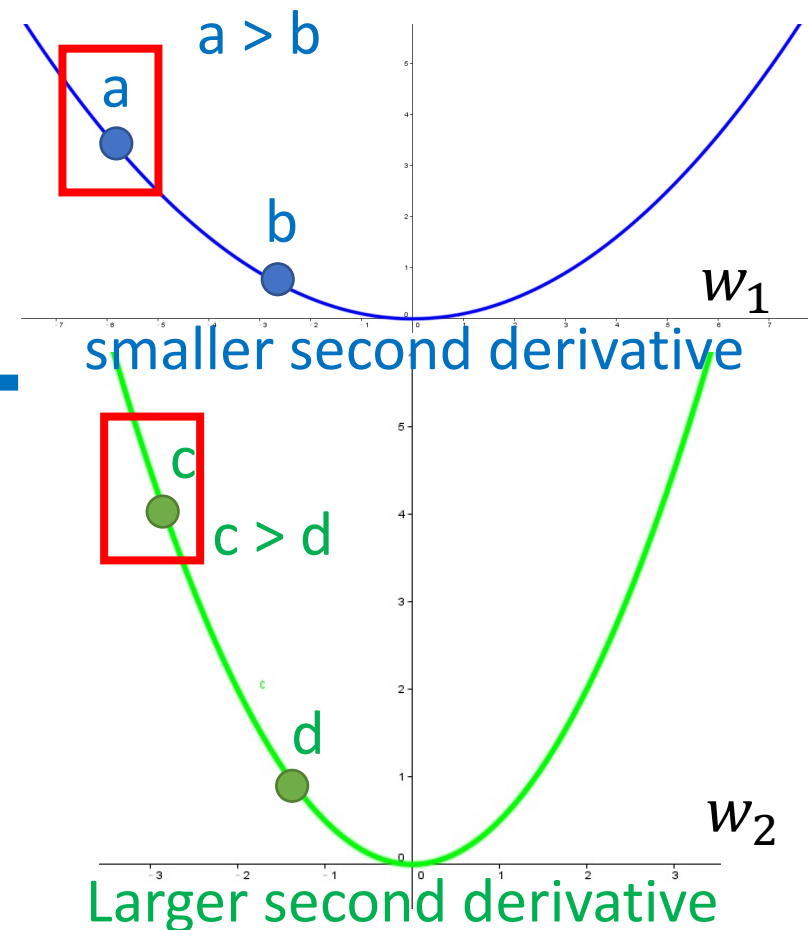
$$\frac{|\text{First derivative}|}{\text{Second derivative}}$$

Comparison between different parameters



~~Larger 1st order derivative means far from the minima~~

Do not cross parameters



$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

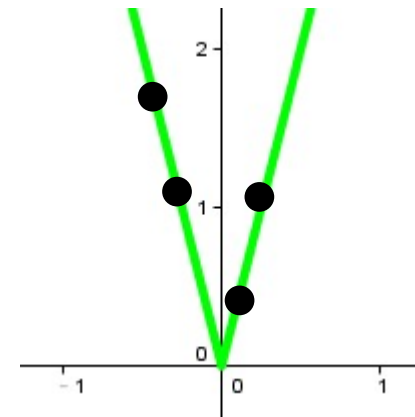
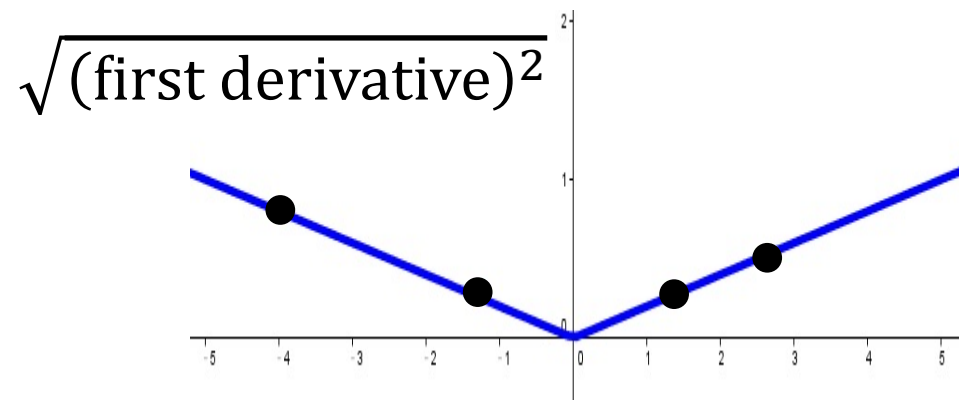
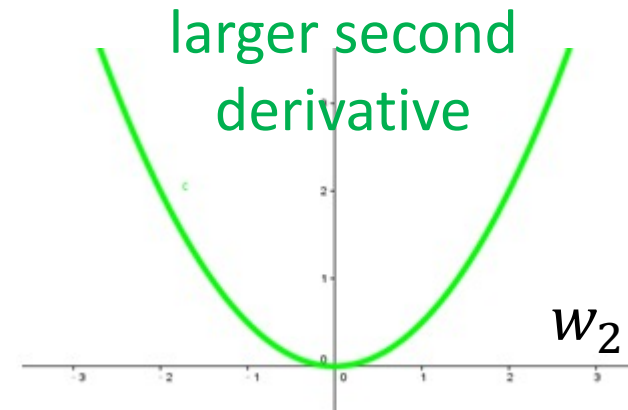
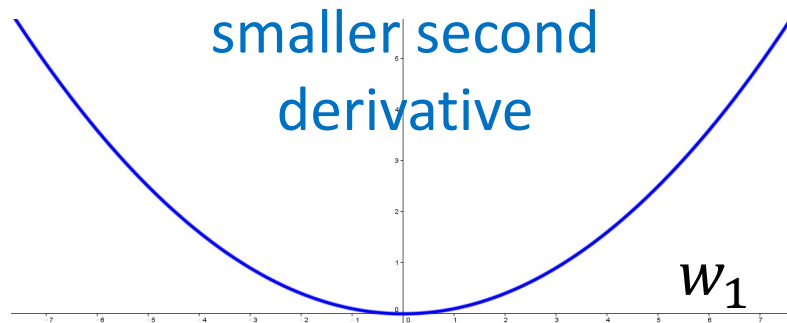
The best step is

|First derivative|

Second derivative

?

Use *first derivative* to estimate *second derivative*



Implement Adagrad by NumPy

```
import numpy as np
```

```
Load features and labels as numpy arrays
```

```
    e.g., X = np.array(X)
```

```
Initial bias (b), weight (w), learning rate (lr), and iteration
```

```
lr_b, lr_w = 0.0, 0.0 // learning rate of bias and weight
```

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

```
    b_grad, w_grad = 0.0, 0.0
```

```
    for j in range(length of X):
```

```
        b_grad = b_grad - 2.0*(y[j] - b - w * X[j]) * 1
```

```
        w_grad = w_grad - 2.0*(y[j] - b - w * X[j]) * X[j]
```

You can also use dot product:
`y = np.dot(x,w) + b`

```
lr_b = lr_b + b_grad ** 2 // customized learning rate
```

```
lr_w = lr_w + w_grad ** 2
```

```
b = b - lr / np.sqrt(lr_b) * b_grad
```

```
w = w - lr / np.sqrt(lr_w) * w_grad
```

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

Gradient Descent

Tip 2: Stochastic Gradient Descent

Make the training faster

Stochastic Gradient Descent

$$L = \sum_n \left(\hat{y}^n - \left(b + \sum w_i x_i^n \right) \right)^2$$

Loss is the summation over all training examples

◆ **Gradient Descent** $\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$

◆ **Stochastic Gradient Descent**

Faster!

Pick an example x^n

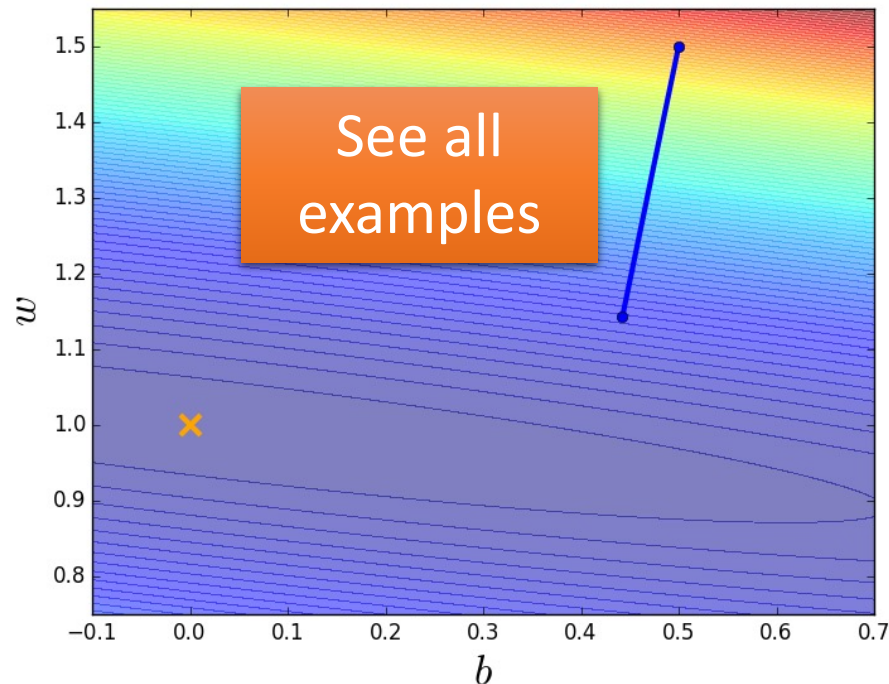
$$L^n = \left(\hat{y}^n - \left(b + \sum w_i x_i^n \right) \right)^2 \quad \theta^i = \theta^{i-1} - \eta \nabla L^n(\theta^{i-1})$$

Loss for only one example

Stochastic Gradient Descent

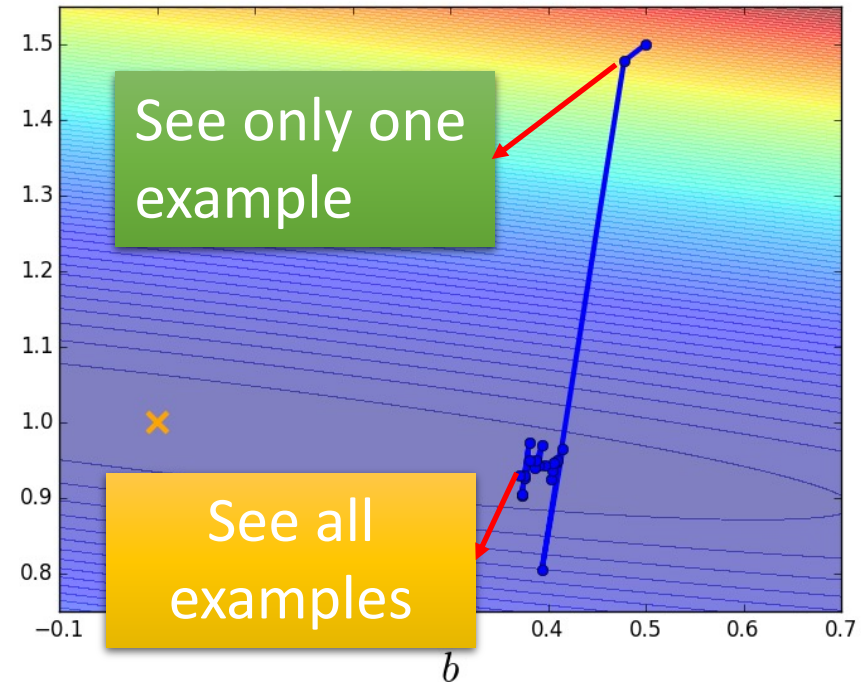
Gradient Descent

Update after seeing all examples



Stochastic Gradient Descent

Update for each example
If there are 20 examples,
20 times faster.



Batch

$$\text{gradient } \mathbf{g} = \begin{bmatrix} \frac{\partial L}{\partial \theta_1} |_{\theta=\theta^0} \\ \frac{\partial L}{\partial \theta_2} |_{\theta=\theta^0} \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \\ \vdots \end{bmatrix} \leftarrow \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \\ \vdots \end{bmatrix} - \begin{bmatrix} \eta \frac{\partial L}{\partial \theta_1} |_{\theta=\theta^0} \\ \eta \frac{\partial L}{\partial \theta_2} |_{\theta=\theta^0} \\ \vdots \end{bmatrix}$$

$$\boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \eta \mathbf{g}$$

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L$$

1. (Randomly) Pick initial values $\boldsymbol{\theta}^0$

2. Compute gradient $\mathbf{g} = \nabla L^1(\boldsymbol{\theta}^0)$

$$\text{update } \boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \eta \mathbf{g}$$

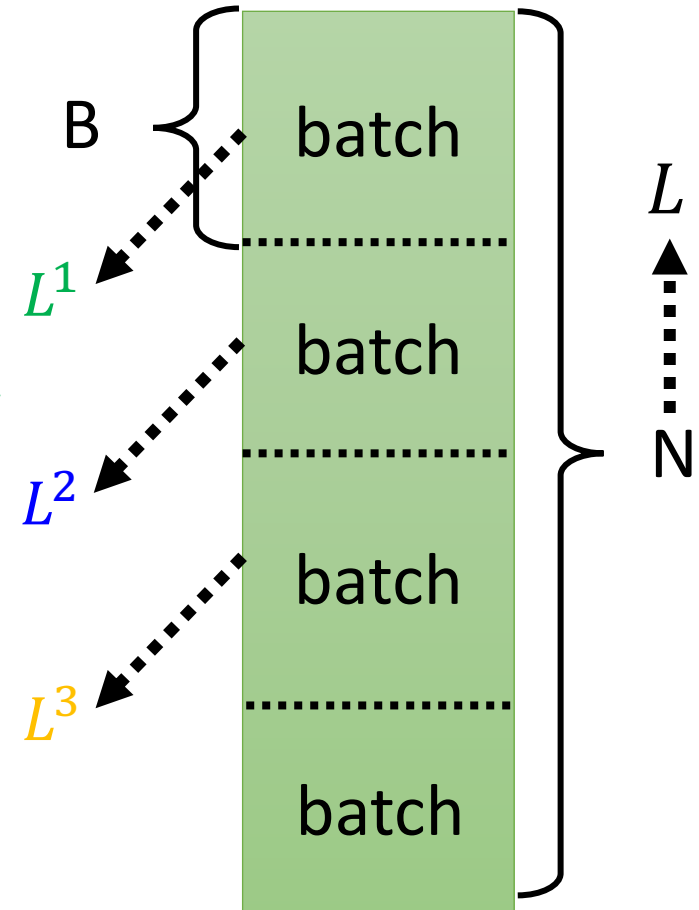
3. Compute gradient $\mathbf{g} = \nabla L^2(\boldsymbol{\theta}^1)$

$$\text{update } \boldsymbol{\theta}^2 \leftarrow \boldsymbol{\theta}^1 - \eta \mathbf{g}$$

4. Compute gradient $\mathbf{g} = \nabla L^3(\boldsymbol{\theta}^2)$

$$\text{update } \boldsymbol{\theta}^3 \leftarrow \boldsymbol{\theta}^2 - \eta \mathbf{g}$$

1 **epoch** = see all the batches once



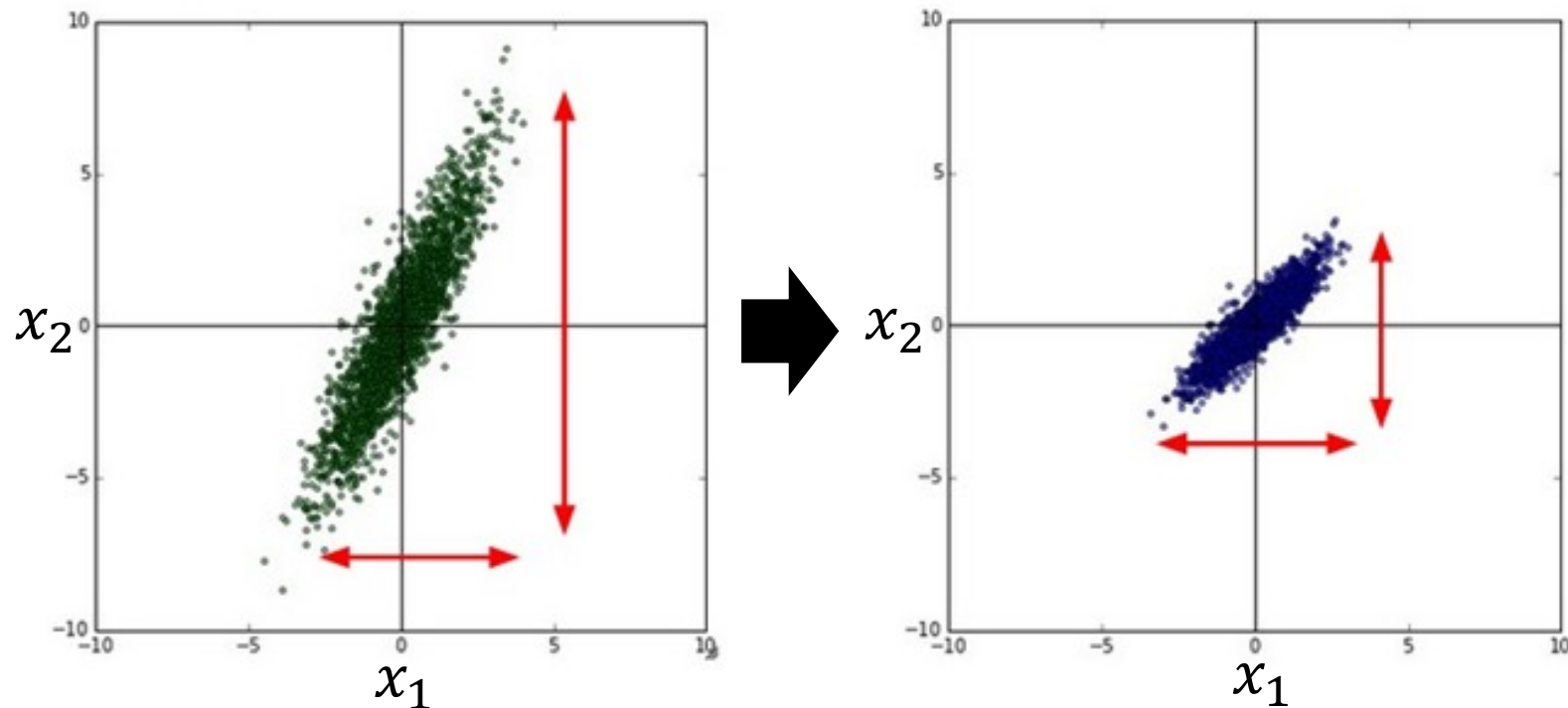
Gradient Descent

Tip 3: Feature Scaling

Feature Scaling

Source of figure:
<http://cs231n.github.io/neural-networks-2/>

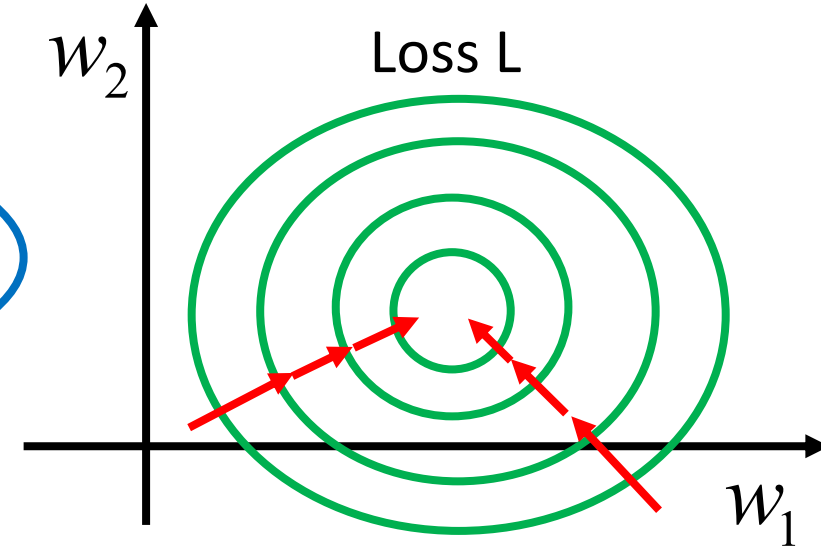
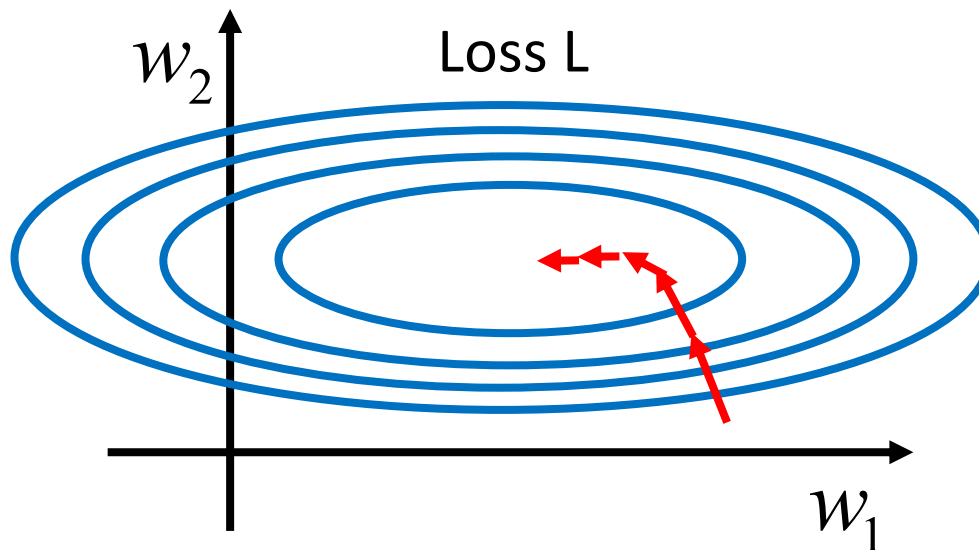
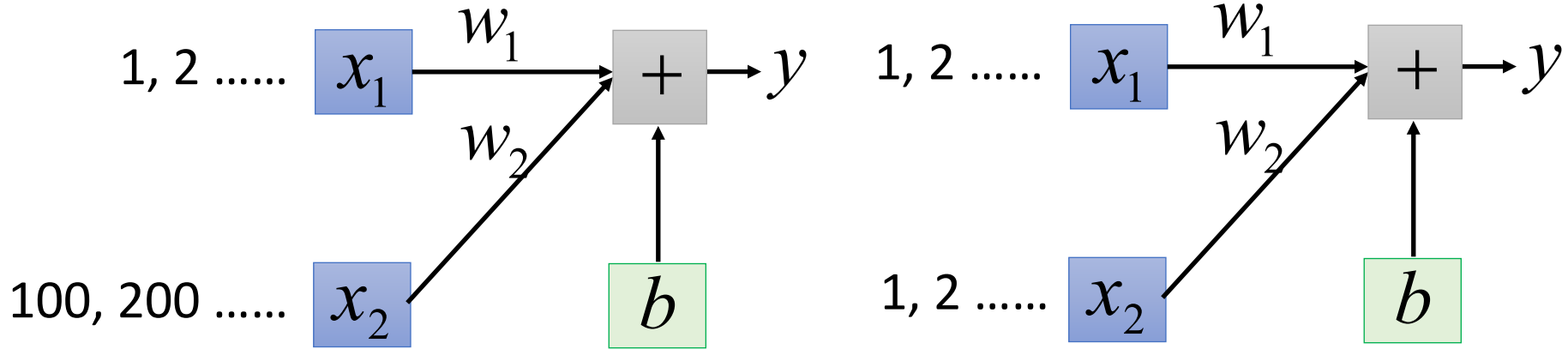
$$y = b + w_1x_1 + w_2x_2$$



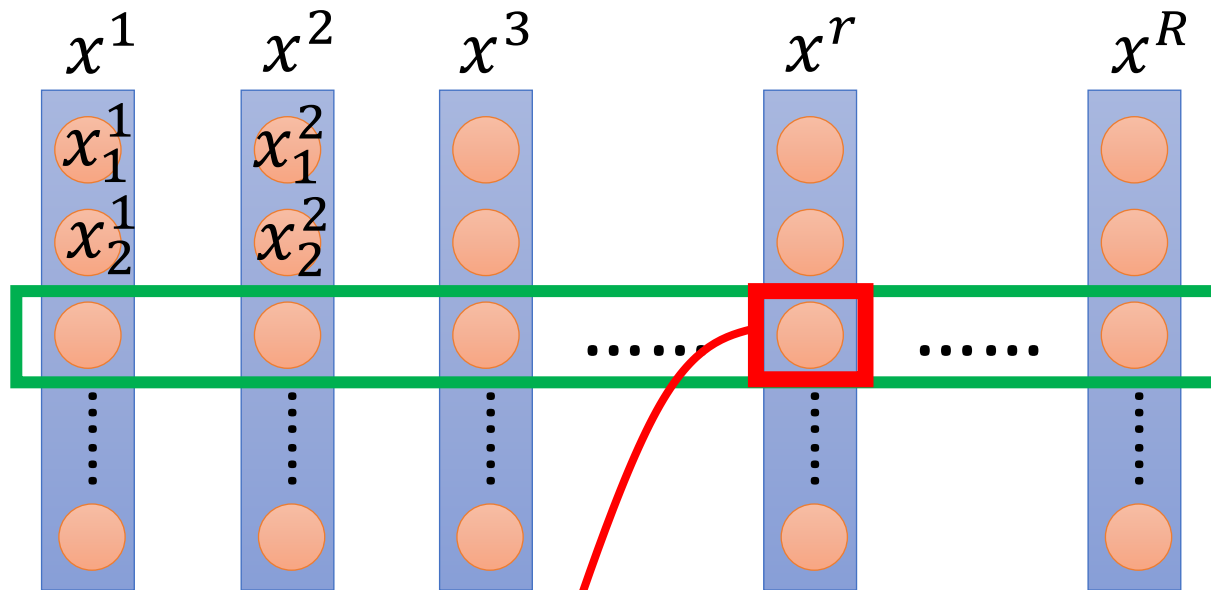
Make different features have the same scaling

Feature Scaling

$$y = b + w_1x_1 + w_2x_2$$



Feature Scaling



For each
dimension i :
mean: m_i
standard
deviation: σ_i

$$x_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$

The means of all dimensions are 0,
and the variances are all 1

Optimization

$$w^*, b^* = \arg \min_{w, b} L(w, b)$$

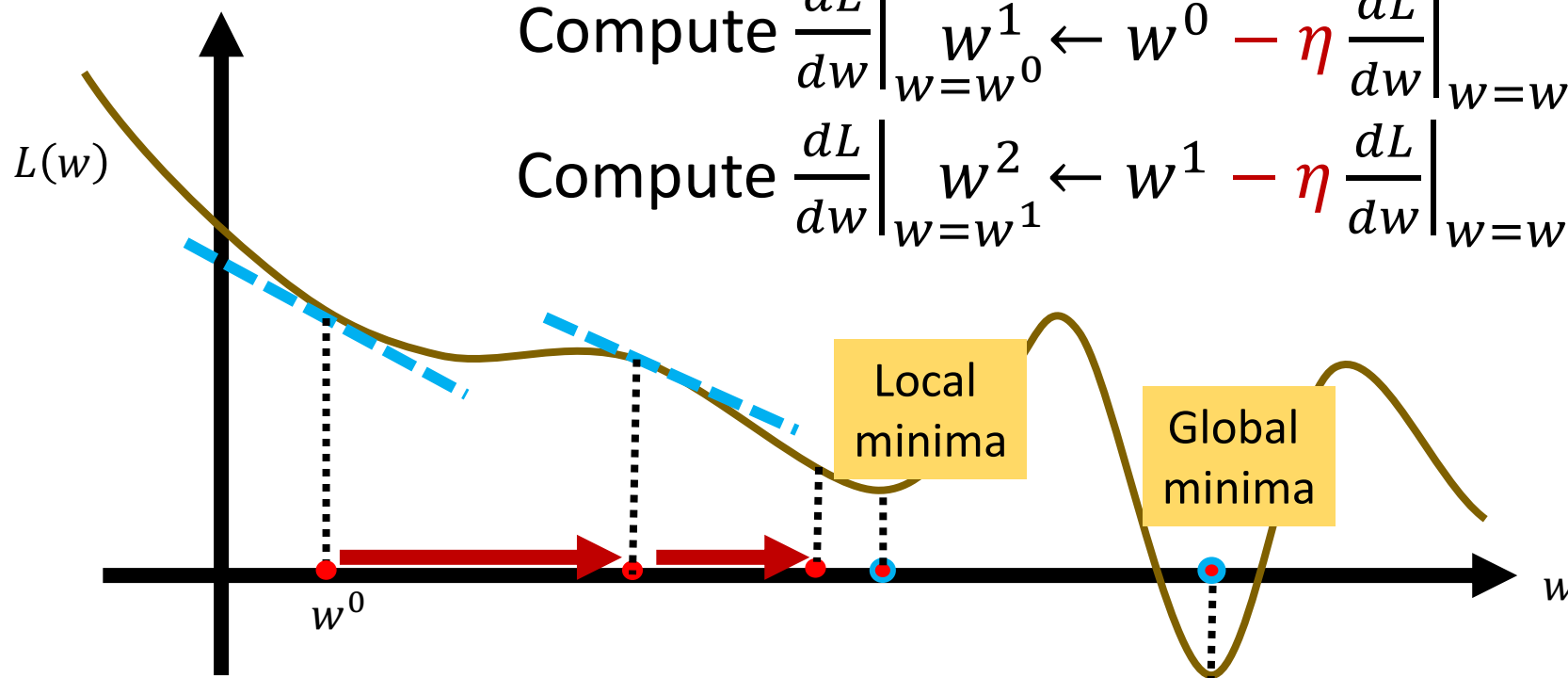
Gradient Descent

Pick an initial value w^0

$$\text{Compute } \frac{dL}{dw} \Big|_{w=w^0} \quad w^1 \leftarrow w^0 - \eta \frac{dL}{dw} \Big|_{w=w^0}$$

$$\text{Compute } \frac{dL}{dw} \Big|_{w=w^1} \quad w^2 \leftarrow w^1 - \eta \frac{dL}{dw} \Big|_{w=w^1}$$

Many
iterations



Reference

- Lecture Slides from Machine Learning by Prof. Hung-Yi Lee
 - https://speech.ee.ntu.edu.tw/~tlkagk/courses/ML_2017/Lecture/Gradient%20Descent.pdf