# Gradient Descent

2020/11/25

# Review: Gradient Descent

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

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

suppose that $\theta$ has two variables $\{\theta_1, \theta_2\}$

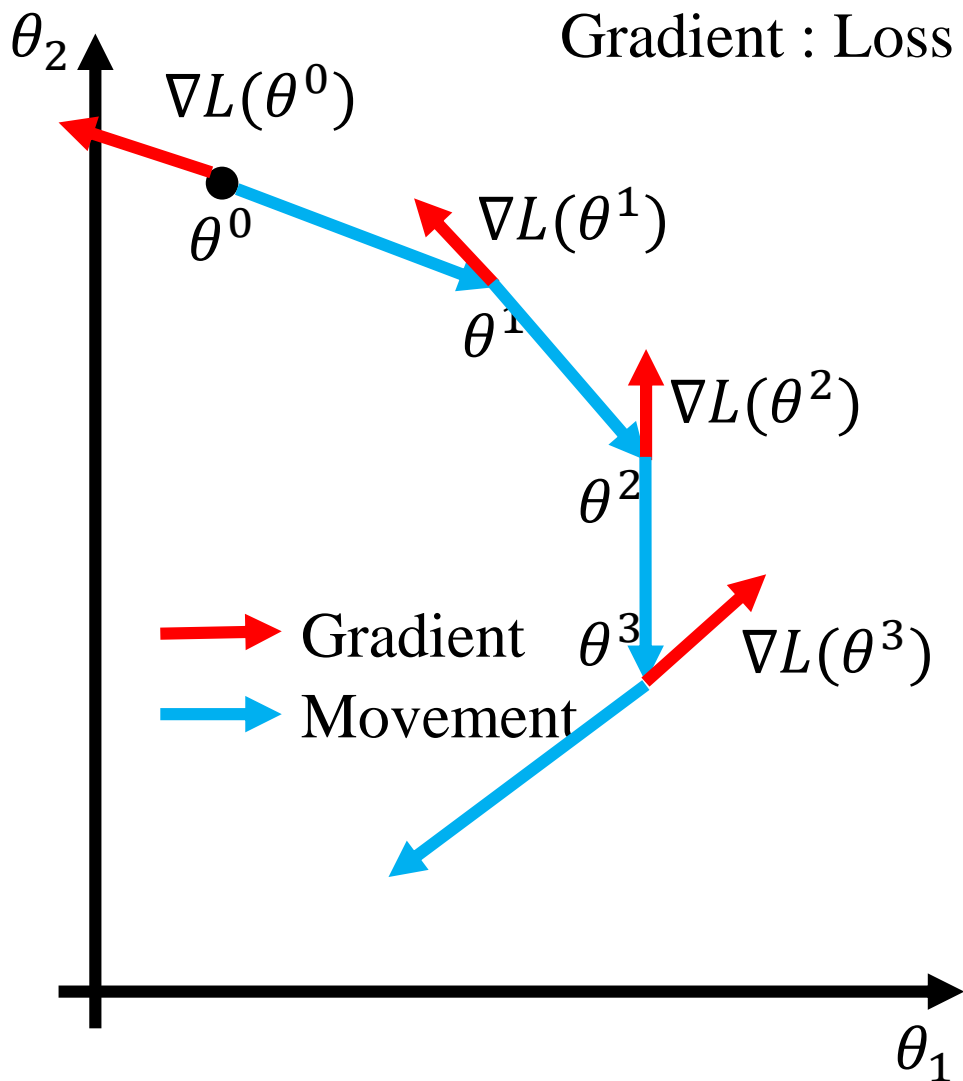Randomly start at $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$

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

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

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

# Review: Gradient Descent

Gradient : Loss 的等高線的法線方向



$\theta_2$

$\nabla L(\theta^0)$

$\theta^0$

$\nabla L(\theta^1)$

$\theta^1$

$\nabla L(\theta^2)$

$\theta^2$

→ Gradient  $\theta^3$  $\nabla L(\theta^3)$

→ Movement

$\theta_1$

Start at position $\theta^0$

Compute gradient at $\theta^0$

Move to $\theta^1 = \theta^0 - \eta \nabla L(\theta^0)$

Compute gradient at $\theta^1$

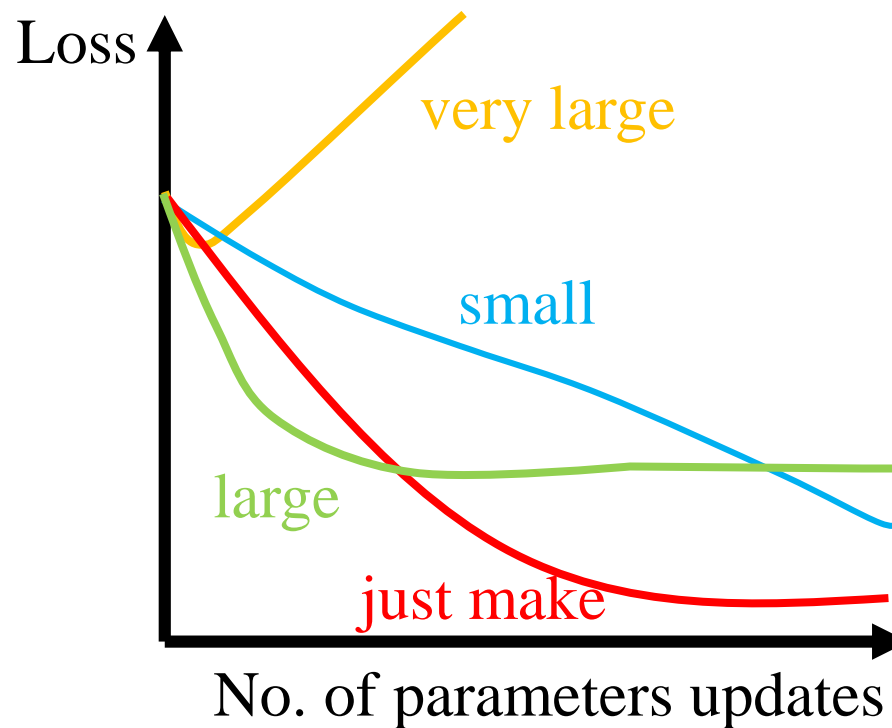Move to $\theta^2 = \theta^1 - \eta \nabla L(\theta^1)$
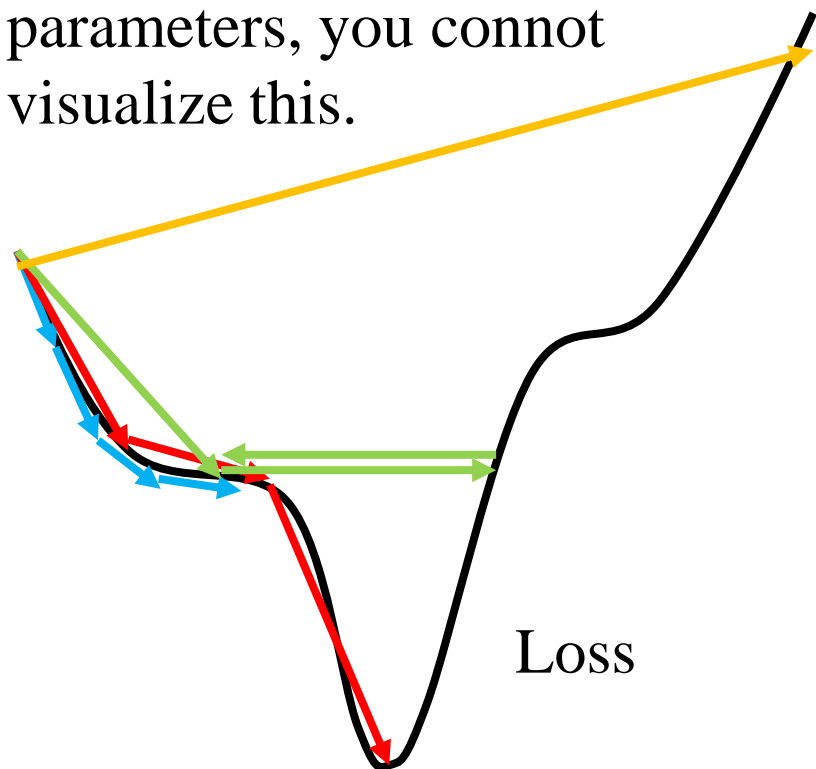
# Gradient Descent
Tip 1 : Tuning your learning rates

# Learning Rate

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

Set the learning rate $\eta$ carefully

If there are more than three parameters, you connot 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*

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

$\sigma^t$: **root mean square** of the previous derivatives of parameter $w$

parameter dependent

# Adagrad

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

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

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

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

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

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

$$\vdots$$

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

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

# Adagrad

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

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

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \qquad \text{1/t decay}$$

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

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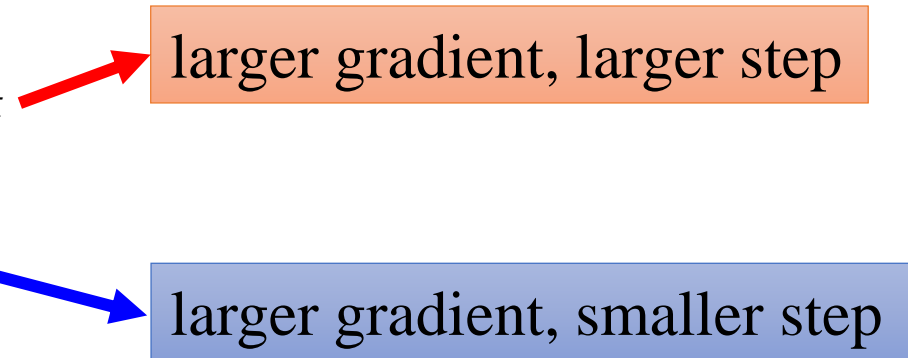
# Contradiction?

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

## *Vanilla Gradient descent*

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

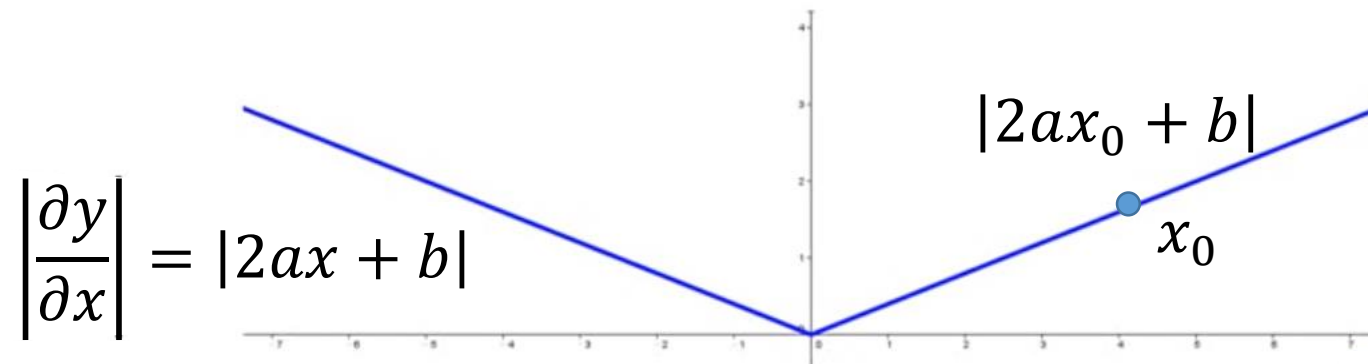→ larger gradient, larger step

## *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. 反差

特别大

| $g^0$ | $g^1$ | $g^2$ | $g^3$ | $g^4$ | ...... |
|-------|-------|-------|-------|-------|--------|
| 0.001 | 0.001 | 0.003 | 0.002 | 0.1 | ...... |

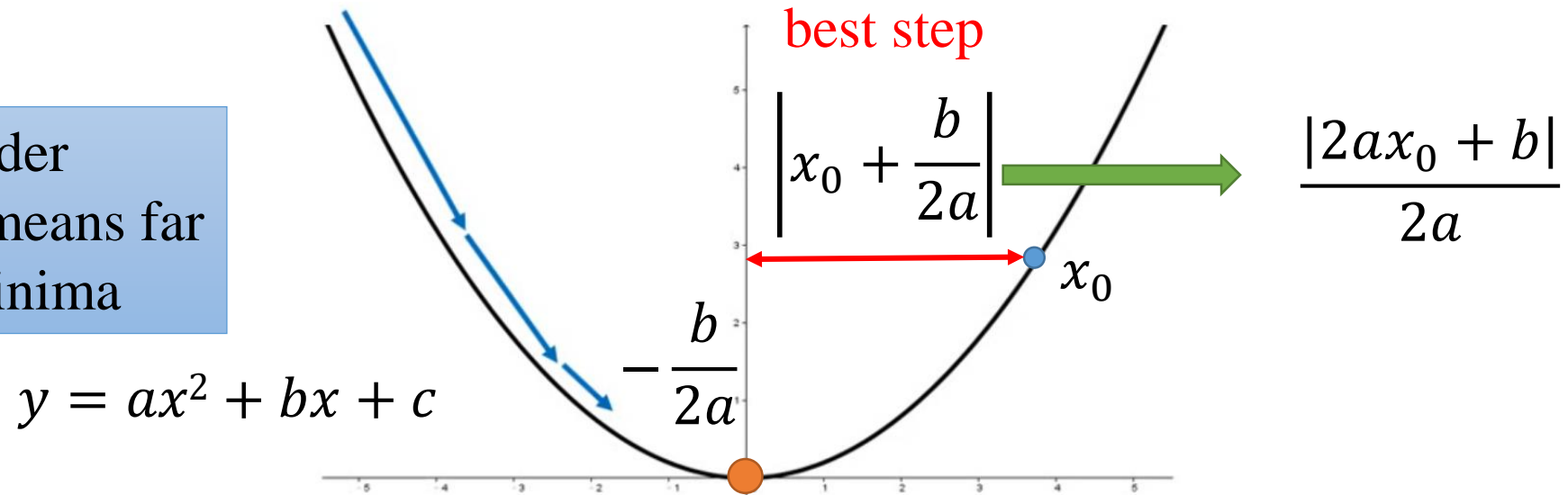| $g^0$ | $g^1$ | $g^2$ | $g^3$ | $g^4$ | ...... |
|-------|-------|-------|-------|-------|--------|
| 10.8 | 20.9 | 31.7 | 12.1 | 0.1 | ...... |

特别小

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^i)^2}} g^t \longrightarrow$$ 造成反差的效果

# larger gradient, larger step?

larger 1st order derivative means far from the minima

$$y = ax^2 + bx + c$$

best step

$$\left| x_0 + \frac{b}{2a} \right|$$

$$\frac{|2ax_0 + b|}{2a}$$

$$-\frac{b}{2a}$$

$$x_0$$

$$\left| \frac{\partial y}{\partial x} \right| = |2ax + b|$$
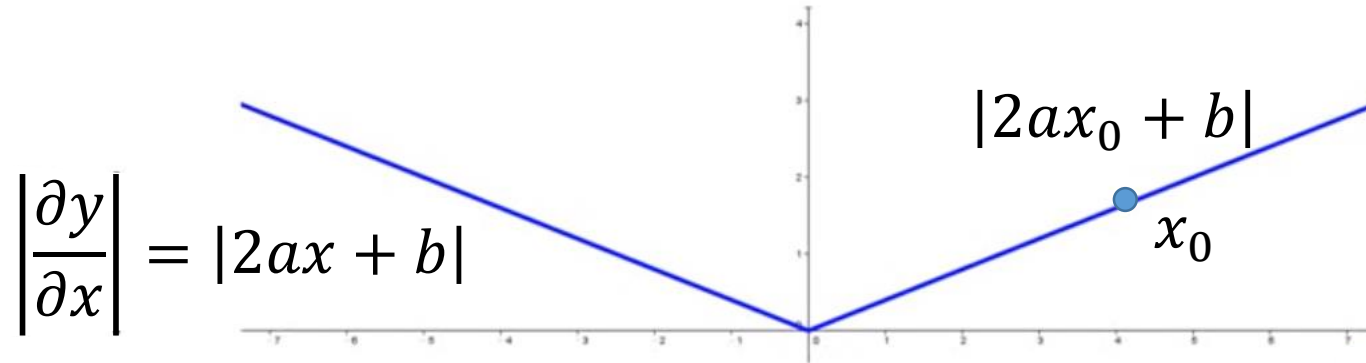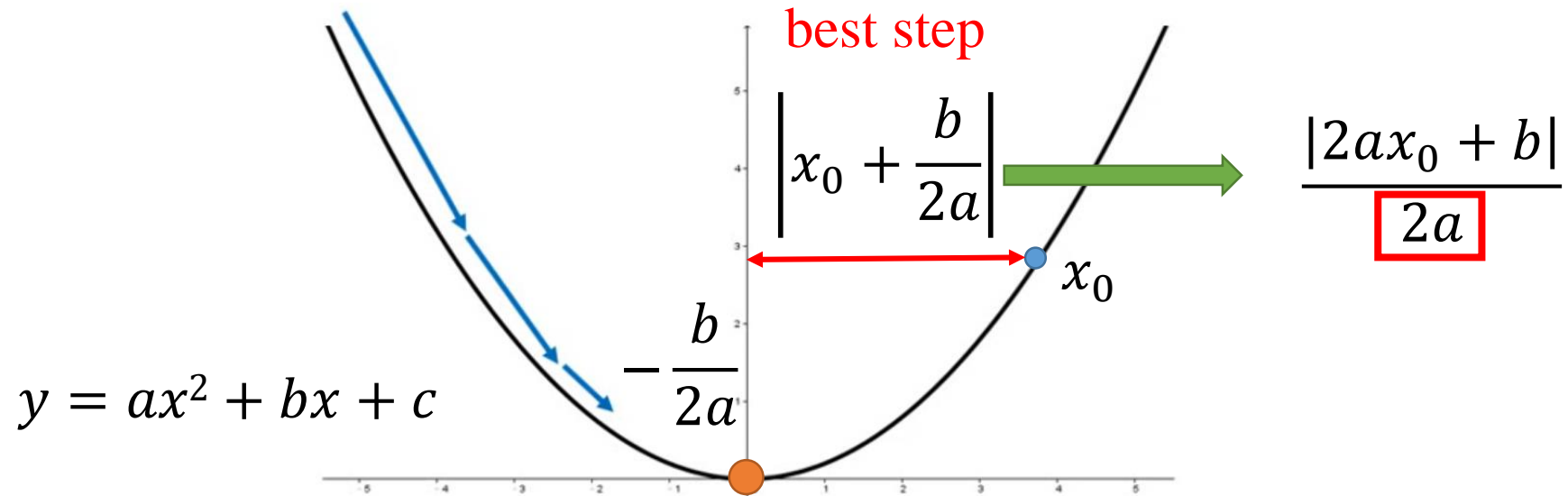
$$|2ax_0 + b|$$

$$x_0$$

# comparison between different parameters

larger 1$^{st}$ order derivative means far from the minima

Do not cross parameters

# second derivative



$$y = ax^2 + bx + c$$

$$\left|\frac{\partial y}{\partial x}\right| = |2ax + b|$$

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

best step

$$\left|x_0 + \frac{b}{2a}\right|$$

$$\frac{|2ax_0 + b|}{\boxed{2a}}$$
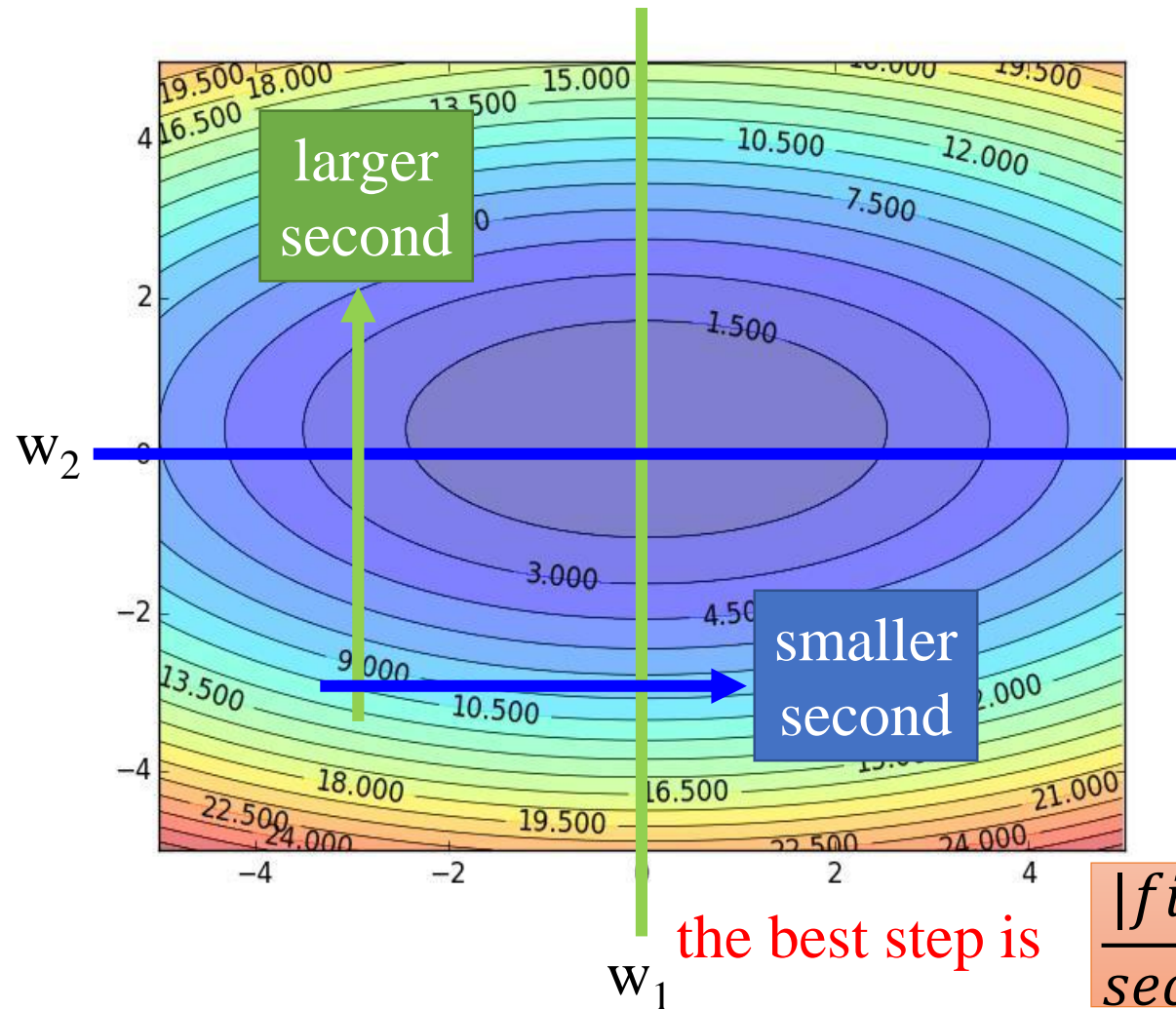
$$-\frac{b}{2a}$$

$x_0$

$$|2ax_0 + b|$$

$x_0$

the best step is $\dfrac{|first\ derivative|}{second\ derivative}$

# comparison between different parameters

larger 1st order derivative means far from the minima

Do not cross parameters

larger second

smaller second

$w_2$

$w_1$

the best step is $\dfrac{|first\ derivative|}{second\ derivative}$

a > b

a

b

$w_1$

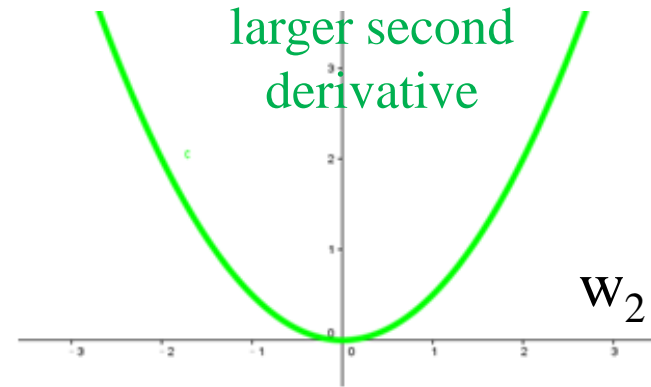smaller second derivative

c > d

c
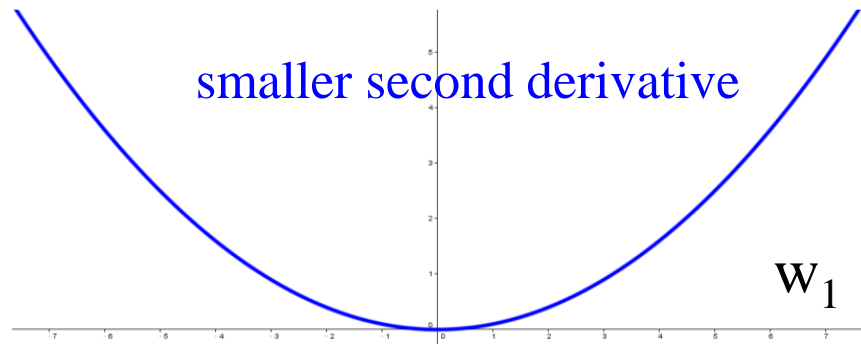
d

larger second derivative

$w_2$

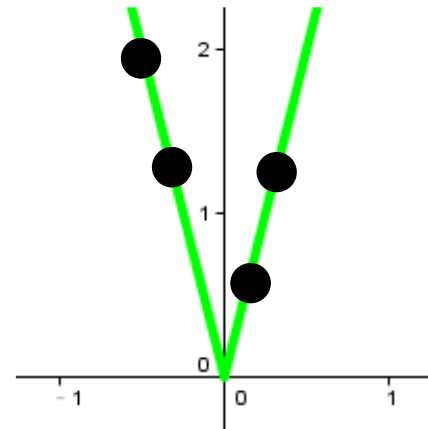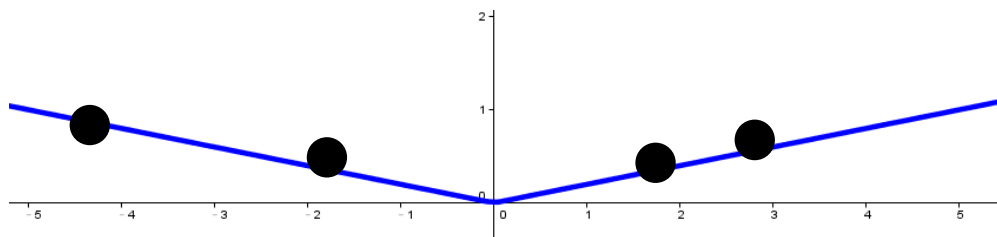$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^i)^2}}g^t$$

the best step is

$$\frac{|first\ derivative|}{second\ derivative}$$

?

Use *first derivative* to estimate *second derivative*



smaller second derivative

$w_1$

larger second derivative

$w_2$

$\sqrt{(first\ dericative)^2}$

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

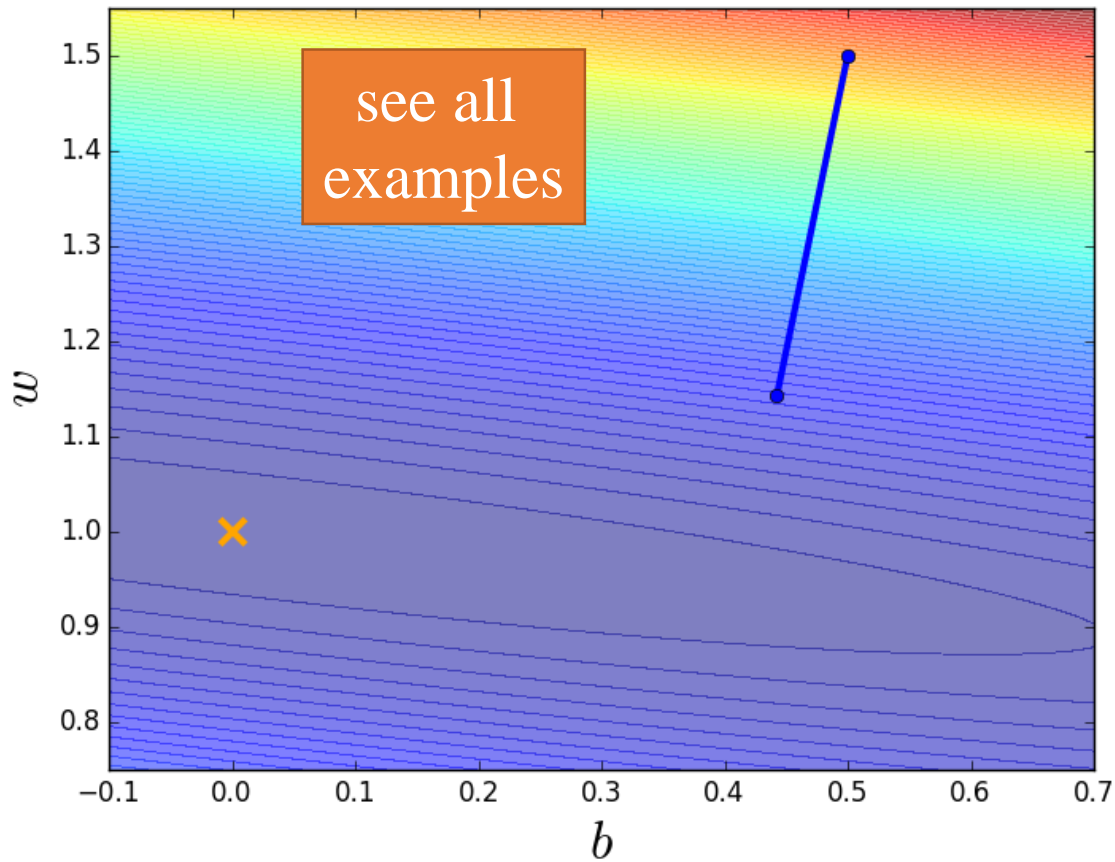◆ **Stochastic Gradient Descent**   Faster!

pick an example $x^n$

$$L^n = \sum_n \left( \hat{y}^n - \left( b + \sum w_i x_i^n \right) \right)^2$$   $\theta^i = \theta^{i-1} - \eta \nabla L^n(\theta^{i-1})$
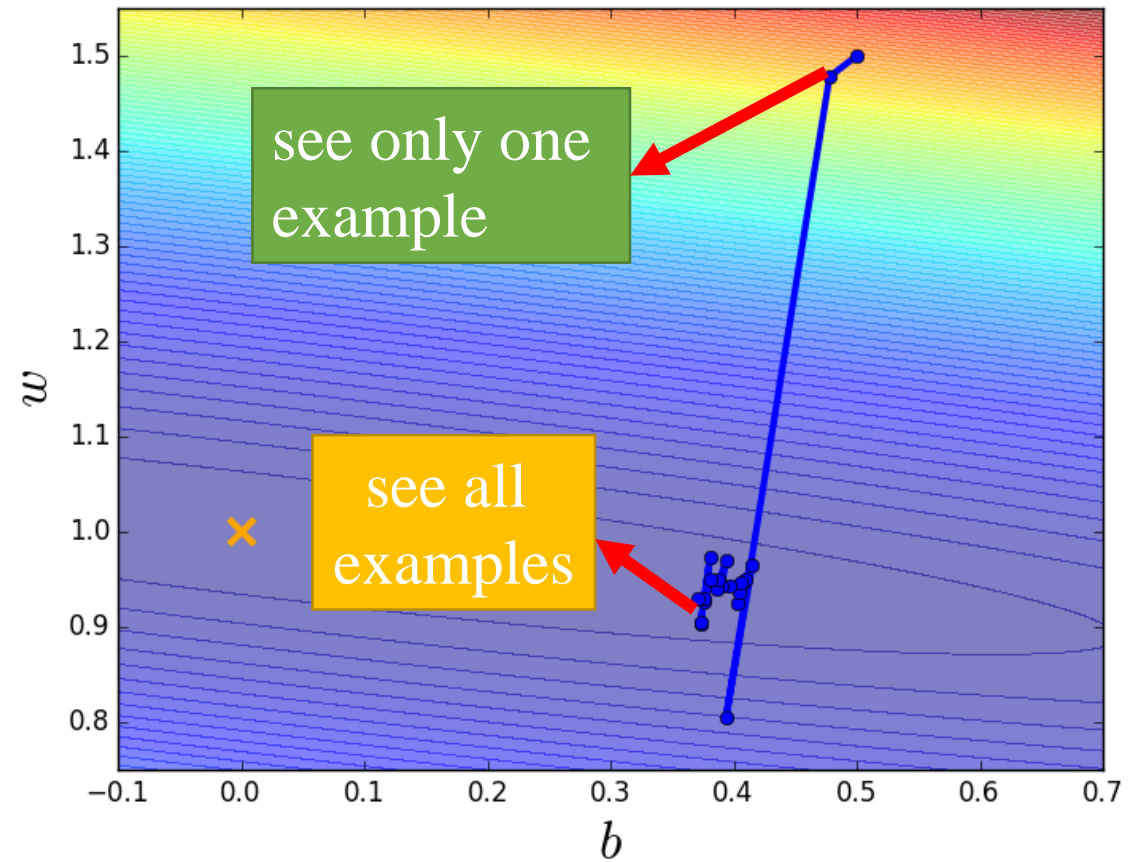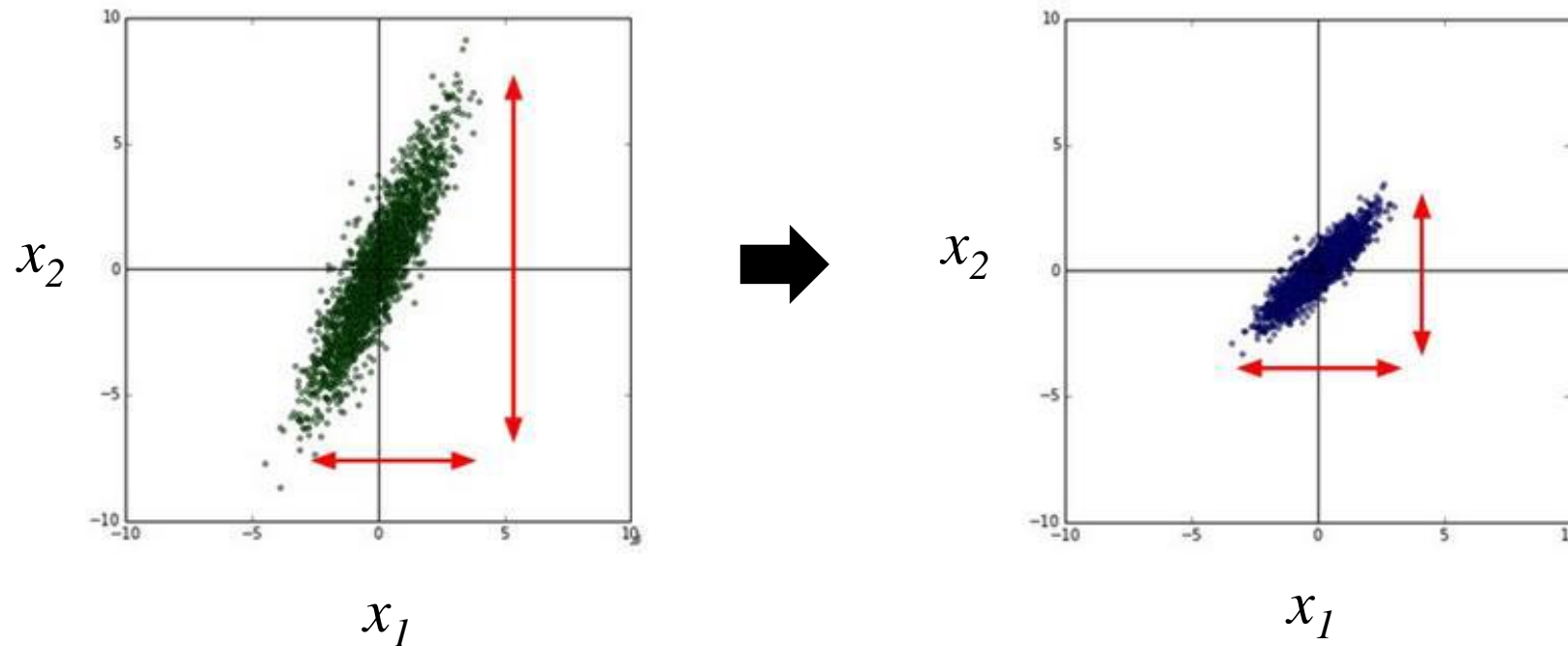
loss for only one example

# Stochastic Gradient Descent

# Gradient Descent
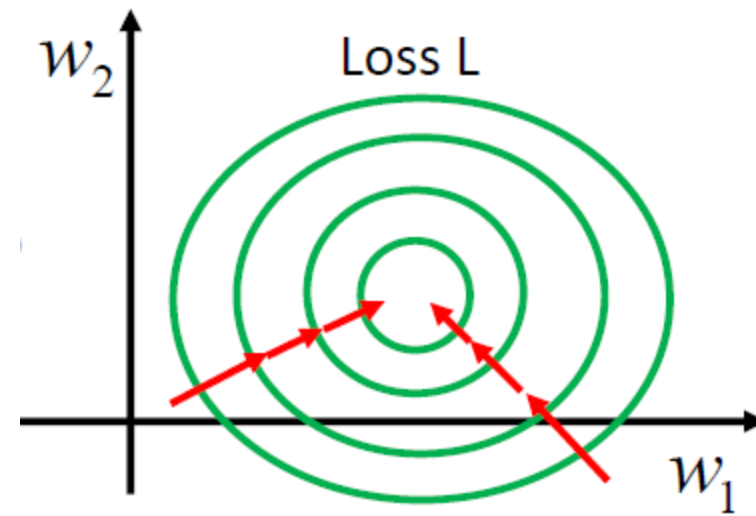# Tip 3: Feature Scaling

# Feature Scaling

$$y = b + w_1 x_1 + w_2 x_2$$



$x_2$     $x_1$      $x_2$     $x_1$

make different features have the same scaling

# Feature Scaling

$$y = b + w_1 x_1 + w_2 x_2$$



1, 2…   $x_1$   $w_1$   $+$ → $y$

100, 200…   $x_2$   $w_2$   $b$

1, 2…   $x_1$   $w_1$   $+$ → $y$

1, 2…   $x_2$   $w_2$   $b$

Loss L

$w_2$   $w_1$

Loss L

$w_2$   $w_1$

# Feature Scaling

$x^1$  $x^2$  $x^3$  $x^r$  $x^R$

for each dimension i:

mean: $m_i$

standard deviation: $\sigma_i$

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

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

# Gradient Descent
## Theory

# Question

- When solving:

$$\theta^* = \arg min_\theta L(\theta) \text{ by gradient descent}$$

- Each time we update the parameters, we obtain $\theta$ that makes $L(\theta)$ smaller.

$$L(\theta^0) > L(\theta^1) > L(\theta^2) > \cdots$$

is this statement correct?

Warning of Math

# Formal Derivation

suppose that $\theta$ has two variable $\{\theta_1, \theta_2\}$

Given a point, we can easily find the point with the smallest value nearby.

how?

# Taylor Series

- **Taylor series**: Let h(x) be any function infinitely differentiable around $x = x_0$.

$$h(x) = \sum_{k=0}^{\infty} \frac{h^{(k)}(x_0)}{k!} (x = x_0)^k$$

$$= h(x_0) + h'(x_0)(x - x_0) + \frac{h''(x_0)}{2!}(x - x_0)^2 + \cdots$$

when x is close to $x_0$ ➡ $h(x) \approx h(x_0) + h'(x_0)(x - x_0)$

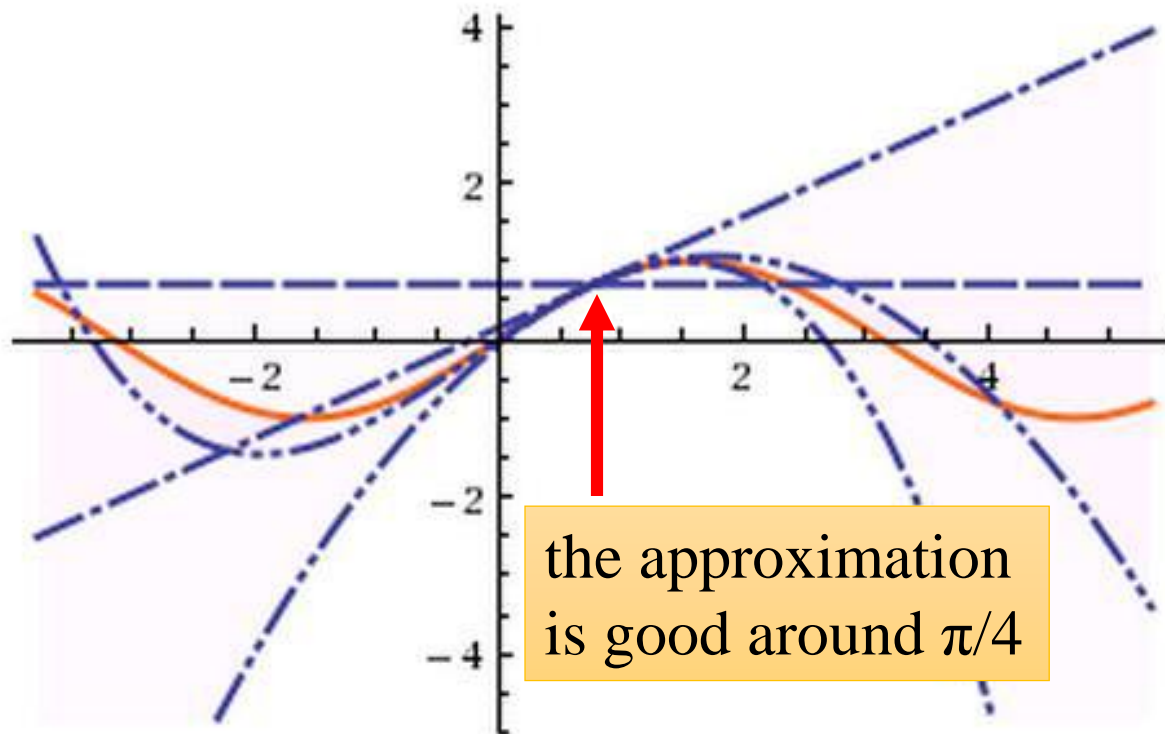e.g. Taylor series for $h(x) = \sin(x)$ around $x_0 = \pi/4$

$$\sin(x) = \frac{1}{\sqrt{2}} + \frac{x - \frac{\pi}{4}}{\sqrt{2}} - \frac{\left(x - \frac{\pi}{4}\right)^2}{2\sqrt{2}} - \frac{\left(x - \frac{\pi}{4}\right)^3}{6\sqrt{2}} + \frac{\left(x - \frac{\pi}{4}\right)^4}{24\sqrt{2}} + \frac{\left(x - \frac{\pi}{4}\right)^5}{120\sqrt{2}} - \frac{\left(x - \frac{\pi}{4}\right)^6}{720\sqrt{2}}$$

$$- \frac{\left(x - \frac{\pi}{4}\right)^7}{5040\sqrt{2}} + \frac{\left(x - \frac{\pi}{4}\right)^8}{40320\sqrt{2}} + \frac{\left(x - \frac{\pi}{4}\right)^9}{362880\sqrt{2}} - \frac{\left(x - \frac{\pi}{4}\right)^{10}}{3628800\sqrt{2}} + \cdots$$



the approximation is good around π/4

# Multivariable Taylor Series

$$h(x, y) = h(x_0, y_0) + \frac{\partial h(x_0, y_0)}{\partial x}(x - x_0) + \frac{\partial h(x_0, y_0)}{\partial y}(y - y_0)$$
$$+ \ something \ related \ to \ (x - x_0)^2 \ and \ (y - y_0)^2 + \cdots$$

When x and y is close to $x_0$ and $y_0$

$$h(x, y) \approx h(x_0, y_0) + \frac{\partial h(x_0, y_0)}{\partial x}(x - x_0) + \frac{\partial h(x_0, y_0)}{\partial y}(y - y_0)$$

# back to formal derivation

based on Taylor Series
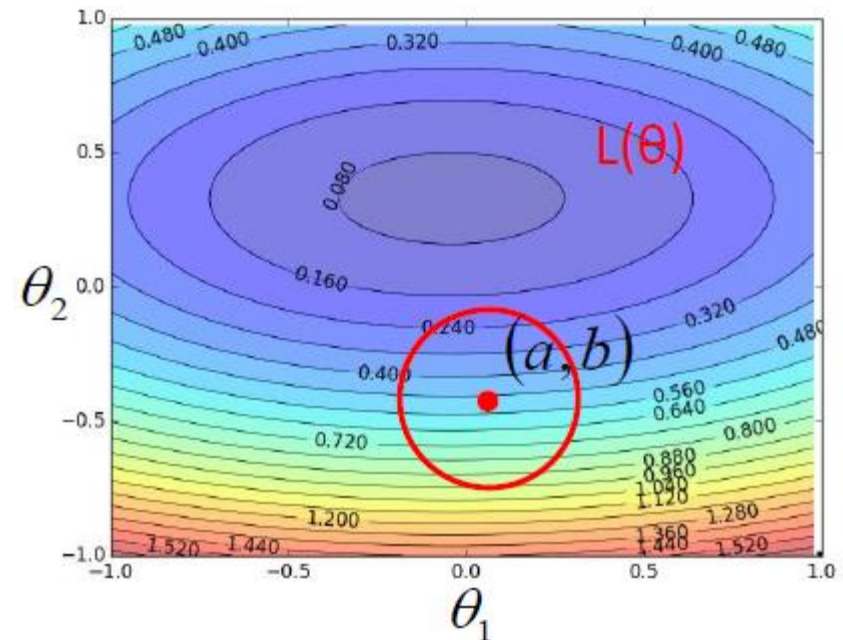
if the red circle is ***small enough***, in the red circle:

$$L(\theta) \approx L(a, b) + \frac{\partial L(a, b)}{\partial \theta_1}(\theta_1 - a) + \frac{\partial L(a, b)}{\partial \theta_2}(\theta_2 - b)$$

$$s = L(a, b)$$

$$u = \frac{\partial L(a, b)}{\partial \theta_1} \qquad v = \frac{\partial L(a, b)}{\partial \theta_2}$$

$$L(\theta) \approx s + u(\theta_1 - a) + v(\theta_2 - b)$$

# back to formal derivation

based on Taylor Series

if the red circle is ***small enough***, in the red circle:

$$L(\theta) \approx s + u(\theta_1 - a) + v(\theta_2 - b)$$

find $\theta_1$ and $\theta_2$ in the red circle **minimizing** $L(\theta)$

$$(\theta_1 - a)^2 + (\theta_2 - b)^2 \leq d^2$$

$$s = L(a, b)$$

$$u = \frac{\partial L(a, b)}{\partial \theta_1} \qquad v = \frac{\partial L(a, b)}{\partial \theta_2}$$

# Gradient descent – two variables

Red circle: (if the radius is small)

$$L(\theta) \approx s + u\underbrace{(\theta_1 - a)}_{\Delta\theta_1} + v\underbrace{(\theta_2 - b)}_{\Delta\theta_2}$$

find $\theta_1$ and $\theta_2$ in the red circle **minimizing** $L(\theta)$

$$\underbrace{(\theta_1 - a)^2}_{\Delta\theta_1} + \underbrace{(\theta_2 - b)^2}_{\Delta\theta_2} \leq d^2$$

to minimize $L(\theta)$

$$\begin{bmatrix} \Delta\theta_1 \\ \Delta\theta_2 \end{bmatrix} = -\eta \begin{bmatrix} u \\ v \end{bmatrix} \implies \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} u \\ v \end{bmatrix}$$



$(\Delta\theta_1, \Delta\theta_2)$

$(\Delta\theta_1, \Delta\theta_2)$

$(u, v)$

# back to formal derivation

based on Taylor Series

if the red circle is ***small enough***, in the red circle:

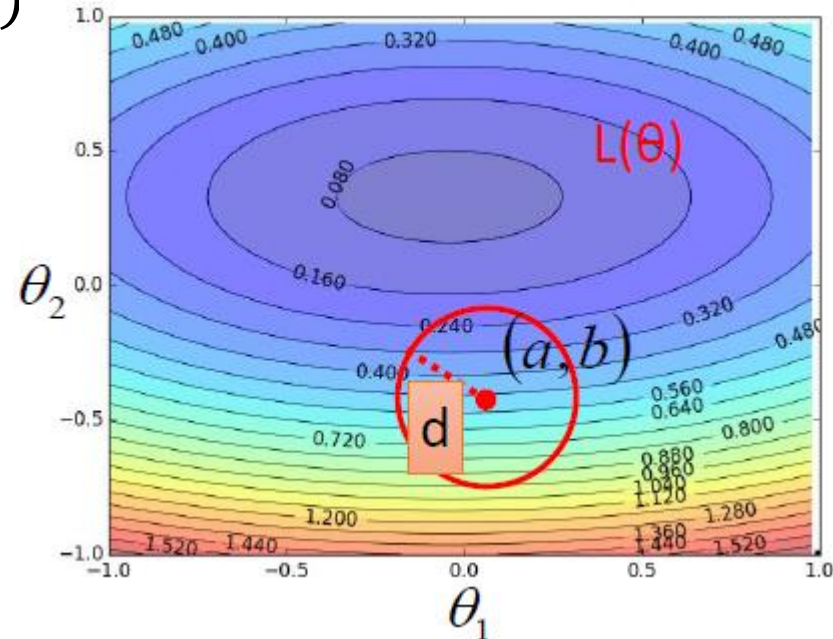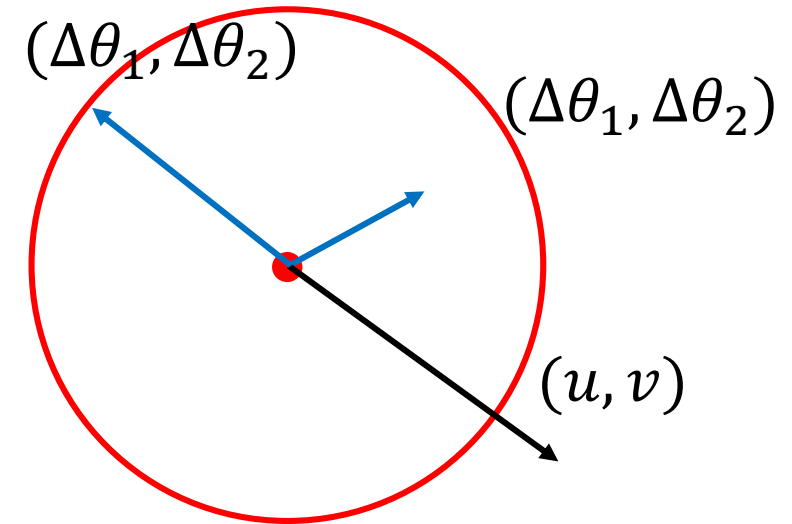$$L(\theta) \approx s + u(\theta_1 - a) + v(\theta_2 - b)$$

$$s = L(a, b)$$

$$u = \frac{\partial L(a, b)}{\partial \theta_1} \qquad v = \frac{\partial L(a, b)}{\partial \theta_2}$$

find $\theta_1$ and $\theta_2$ yielding the smallest value of $L(\theta)$ in the circle

$$\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} \dfrac{\partial L(a, b)}{\partial \theta_1} \\ \dfrac{\partial L(a, b)}{\partial \theta_2} \end{bmatrix}$$ this is gradient descent

not satisfied if the red circle (learning rate) is not small enough

you can consider the second order term, e.g. Newton's method

# More Limitation of Gradient Descent