

Gradient Descent

Review: Gradient Descent

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

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

loss 求最小

Suppose that θ has two variables $\{\theta_1, \theta_2\}$ *写得更简洁.*

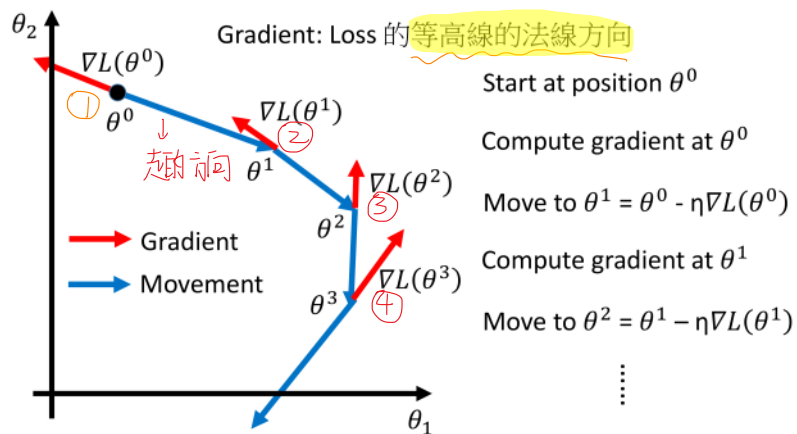
Randomly start at $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$ *就叫 Gradient* $\Rightarrow \nabla L(\theta) = \begin{bmatrix} \partial L(\theta_1)/\partial \theta_1 \\ \partial L(\theta_2)/\partial \theta_2 \end{bmatrix}$

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

L在 θ^0 的参数

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

Review: Gradient Descent



Gradient Descent

Tip 1: Tuning your learning rates

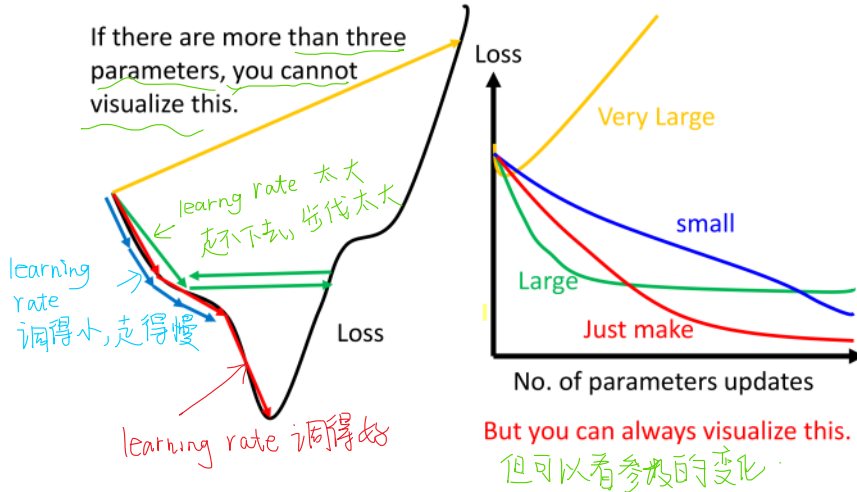
小心

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.



自动调 Learning Rate

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

好几次后

- After several epochs, we are close to the destination, so we reduce the learning rate

- E.g. $1/t$ decay: $\eta^t = \eta / \sqrt{t + 1}$ η 与次数有关.

- Learning rate cannot be one-size-fits-all

- Giving different parameters different learning rates

不同参数给不同的 Learning Rate

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

之前 个微分

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

σ^t

parameter w

Parameter dependent

每个参数的 σ^t 都不同

Adagrad

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

过去算过的微分值的

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

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

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

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

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

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

\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**

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

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

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

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

把 learning rate 直接写成这样

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

取决于 η 和 gradient

Vanilla Gradient descent

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

Larger gradient, larger step

参数更新

Adagrad

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

Larger gradient, larger step

2x17

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

取决于 η 和 gradient

Vanilla Gradient descent

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

Larger gradient, larger step 参数更新

Adagrad

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

Larger gradient, larger step

Larger gradient, smaller step

和我们期望相反

Intuitive Reason $\eta^t = \frac{\eta}{\sqrt{t+1}}$ $g^t = \frac{\partial C(\theta^t)}{\partial w}$

- How surprise it is 反差

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

特别大

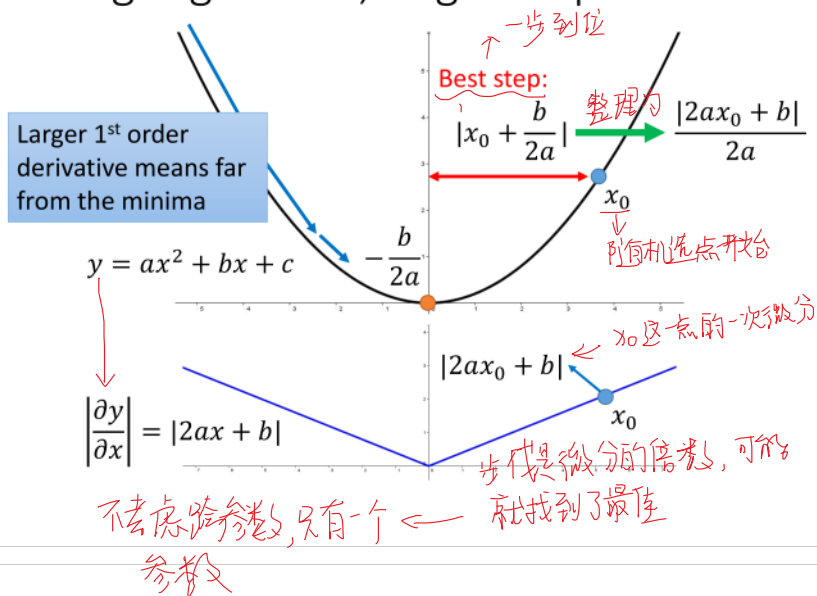
特别小

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

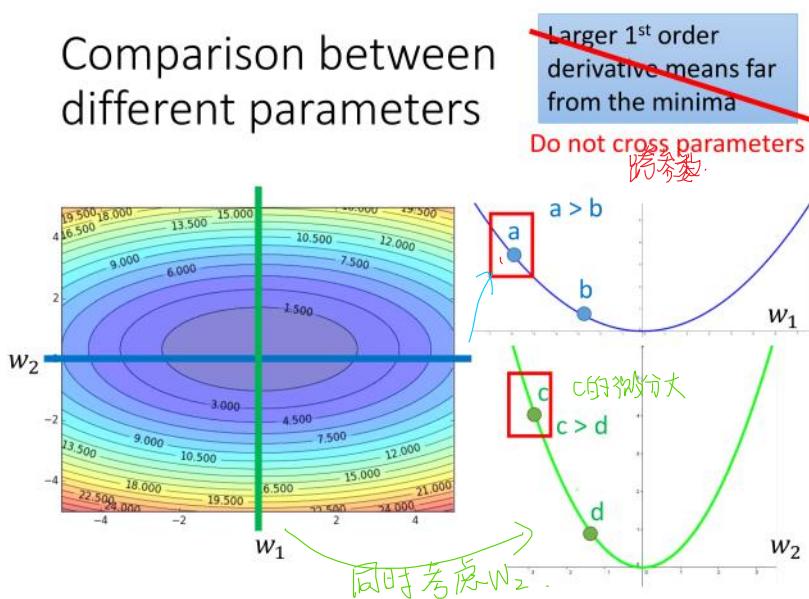
造成反差的效果 多少

过去的 gradient 有多大

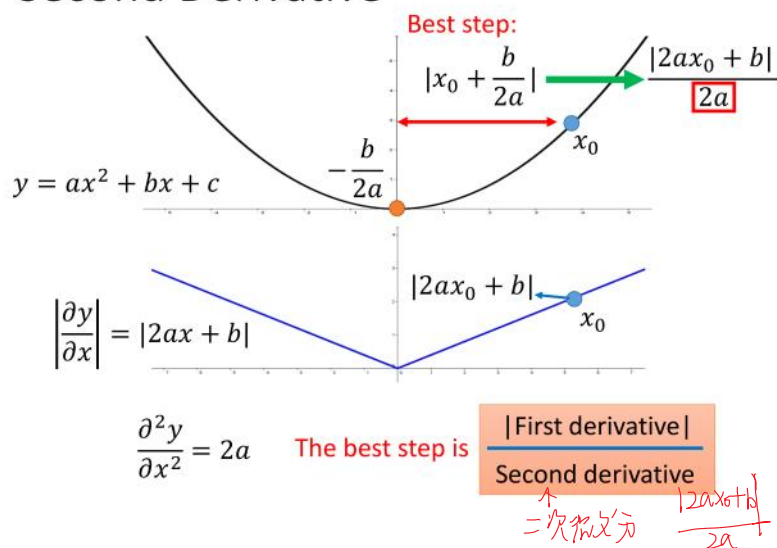
Larger gradient, larger steps?



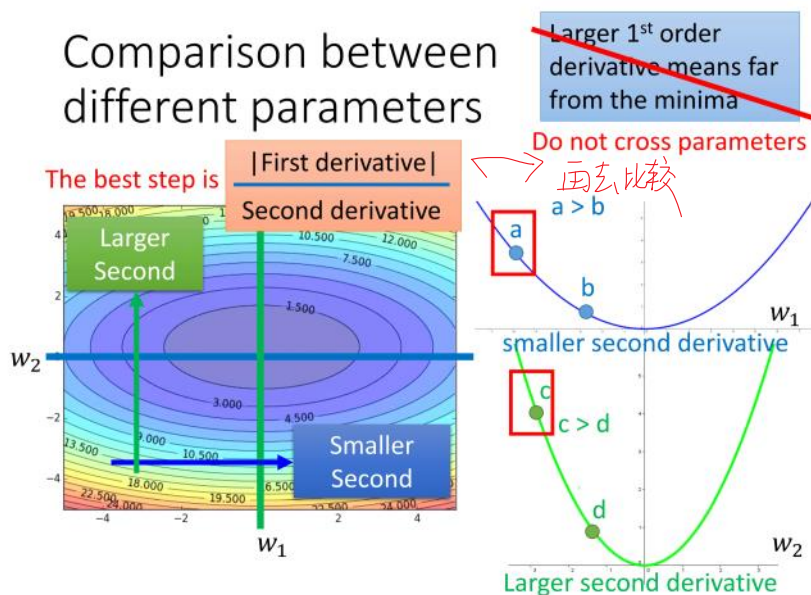
Comparison between different parameters



Second Derivative



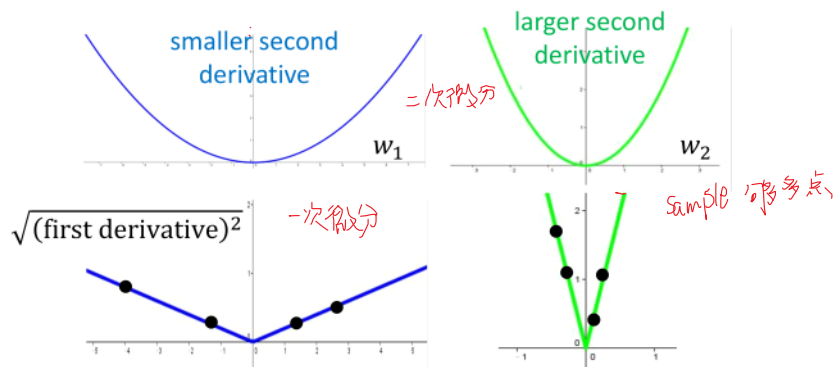
Comparison between different parameters



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

η 常数.
 g^t First derivative
 $\sqrt{\sum_{i=0}^t (g^i)^2}$ Second derivative
 ?
 越累代表二次微分

Use first derivative to estimate second derivative



Gradient Descent

Tip 2: Stochastic Gradient Descent

Make the training faster

让 training 更快

Stochastic Gradient Descent

Regression 的 Loss ↓

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

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

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

Loss for only one example

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

每次拿一个 example 的 Loss

每次拿一个 example 出来 \triangle 随机取
不用累加了 \leftarrow 顺序取

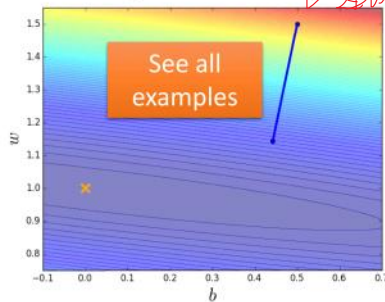
- Demo

Stochastic Gradient Descent

^{SGD} Stochastic Gradient Descent

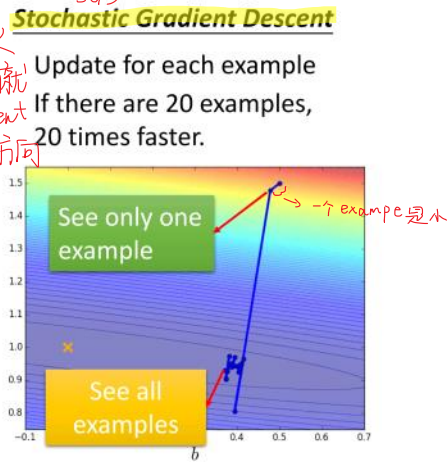
Gradient Descent

Update after seeing all examples



Update for each example

If there are 20 examples, 20 times faster.



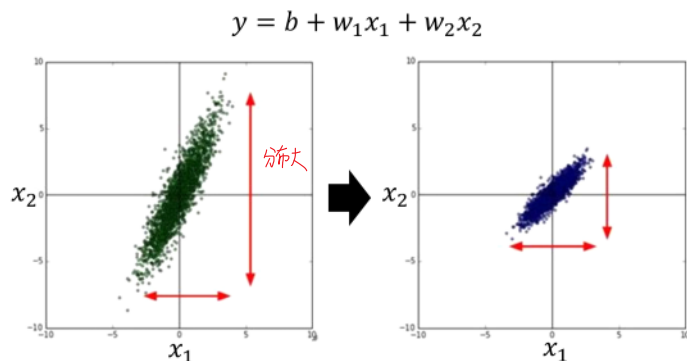
Gradient Descent

Tip 3: Feature Scaling

特征缩放
特征归一化

Feature Scaling

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



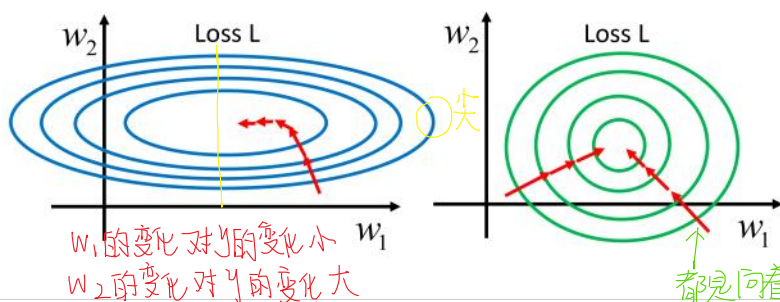
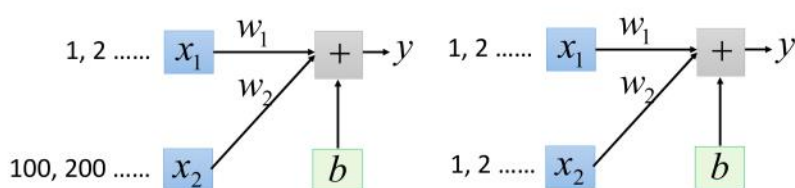
Make different features have the same scaling

把他们的range 放到一样

为什么要做

Feature Scaling 呢?

$$y = b + w_1x_1 + w_2x_2$$

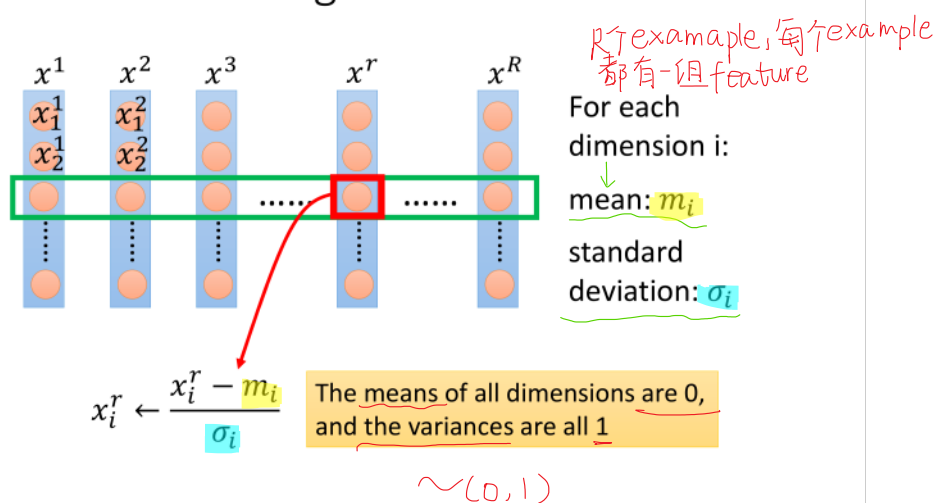


w_1 的变化对 y 的变化小
 w_2 的变化对 y 的变化大

w_1 对 y 有差不多的影响
 w_2 对 y 有差不多的影响

都是向着圆心走

Feature Scaling



Gradient Descent

Theory

Handwritten note: 理论基础

Question

- When solving:

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

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

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

Is this statement correct?

这个陈述是对的吗?

不对

更新后, Loss 不一定下降

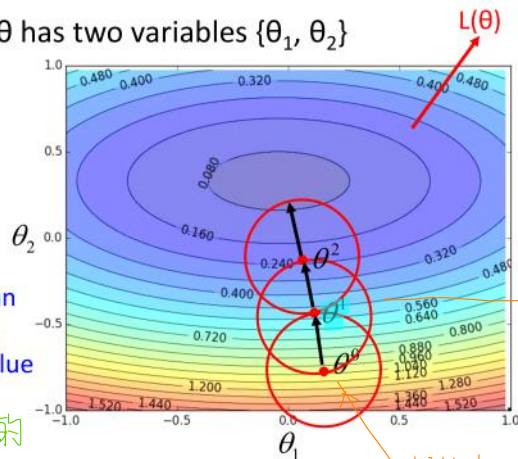
Warning of Math

Formal Derivation

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

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

如何快速找到最小的那个点呢?



再画圈, 再在圆圈内找最低点 θ^2 .

给起始点 θ^0 , 在红圈內找一个最低点 θ^1

Taylor Series

泰勒展开 泰勒级数

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

把 $h(x)$ 写成

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

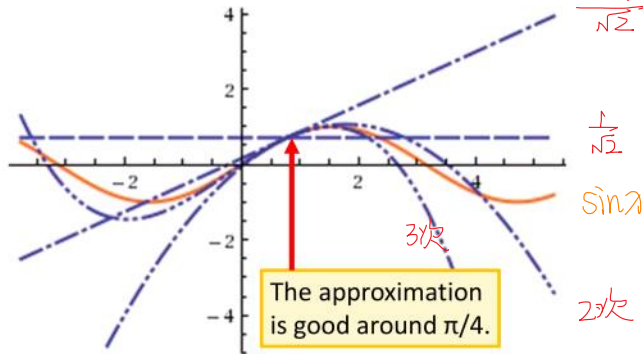
$$= h(x_0) + h'(x_0)(x - x_0) + \frac{h''(x_0)}{2!} (x - x_0)^2 + \dots$$

When x is close to x_0 \Rightarrow $h(x) \approx h(x_0) + h'(x_0)(x - x_0)$
 $x \rightarrow x_0$

例

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

$$\sin(x) = \frac{1}{\sqrt{2}} + \frac{x - \frac{\pi}{4}}{\sqrt{2}} - \frac{(x - \frac{\pi}{4})^2}{2\sqrt{2}} - \frac{(x - \frac{\pi}{4})^3}{6\sqrt{2}} + \frac{(x - \frac{\pi}{4})^4}{24\sqrt{2}} + \frac{(x - \frac{\pi}{4})^5}{120\sqrt{2}} - \frac{(x - \frac{\pi}{4})^6}{720\sqrt{2}} - \frac{(x - \frac{\pi}{4})^7}{5040\sqrt{2}} + \frac{(x - \frac{\pi}{4})^8}{40320\sqrt{2}} + \frac{(x - \frac{\pi}{4})^9}{362880\sqrt{2}} - \frac{(x - \frac{\pi}{4})^{10}}{3628800\sqrt{2}} + \dots$$



$x \rightarrow \frac{\pi}{4}$ 时, 后面高次项都趋近于 0

多呢

Multivariable Taylor Series

$(x, y) \rightarrow (x_0, y_0)$

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

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 红圈足够小

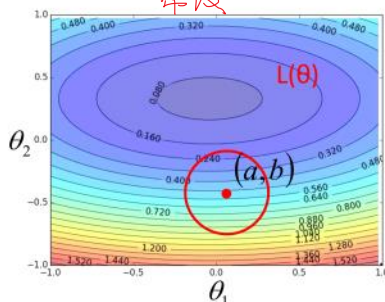
$$L(\theta) \approx \underbrace{L(a,b)}_{\text{常数}} + \underbrace{\frac{\partial L(a,b)}{\partial \theta_1}}_{\text{常数}}(\theta_1 - a) + \underbrace{\frac{\partial L(a,b)}{\partial \theta_2}}_{\text{常数}}(\theta_2 - b) \leftarrow \text{红圈}$$

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

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

$$L(\theta)$$

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



Back to Formal Derivation

Based on Taylor Series:

If the red circle is **small enough**, in the red circle

constant 常数

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

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

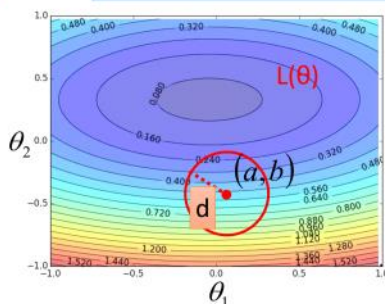
$$u = \frac{\partial L(a,b)}{\partial \theta_1}, v = \frac{\partial L(a,b)}{\partial \theta_2}$$

Find θ_1 and θ_2 in the red circle

minimizing $L(\theta)$

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

Simple, right?



Gradient descent – two variables

Red Circle: (If the radius is small)

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

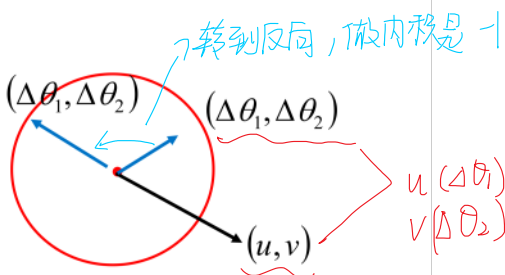
Find θ_1 and θ_2 in the red circle
minimizing $L(\theta)$

$$\left(\frac{\theta_1 - a}{\Delta \theta_1} \right)^2 + \left(\frac{\theta_2 - b}{\Delta \theta_2} \right)^2 \leq d^2$$

To minimize $L(\theta)$

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

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



Back to Formal Derivation

Based on Taylor Series:

If the red circle is **small enough**, in the red circle

constant

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

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

前题该式成立

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

Find θ_1 and θ_2 yielding the smallest value of $L(\theta)$ in the circle

$$\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial L(a, b)}{\partial \theta_1} \\ \frac{\partial L(a, b)}{\partial \theta_2} \end{bmatrix}$$

This is gradient descent.

代进去就是 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**.

二次项

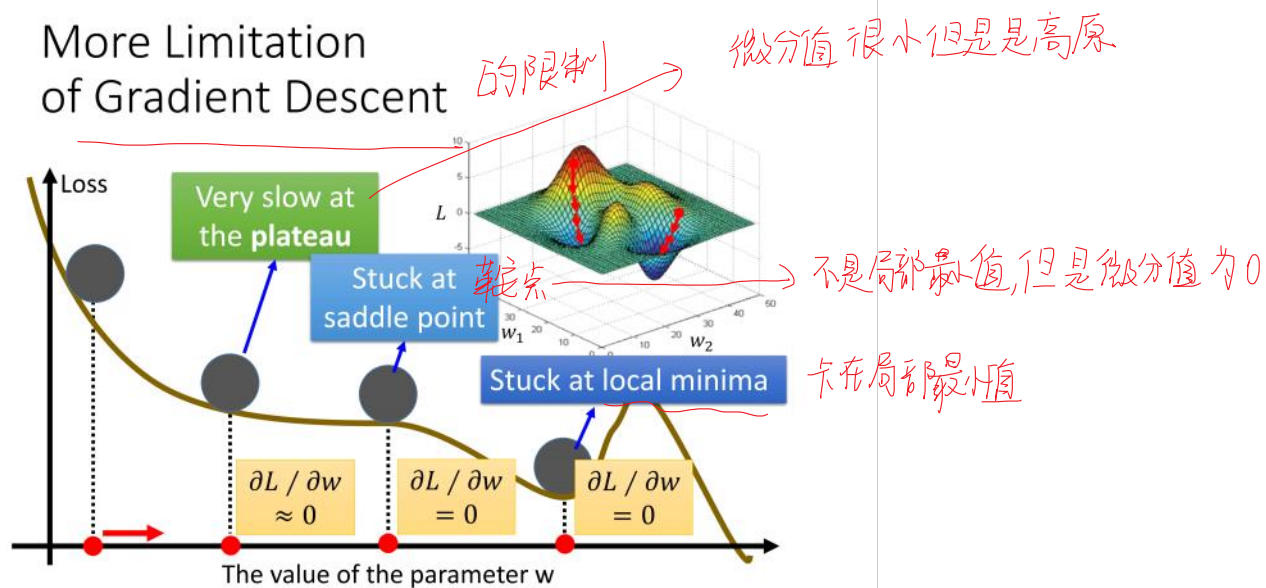
牛顿法

二次微分

learning rate 要设好 → 使 circle 足够小.
使上面的式子成立

End of Warning

More Limitation of Gradient Descent



Acknowledgement

- 感謝 Victor Chen 發現投影片上的打字錯誤