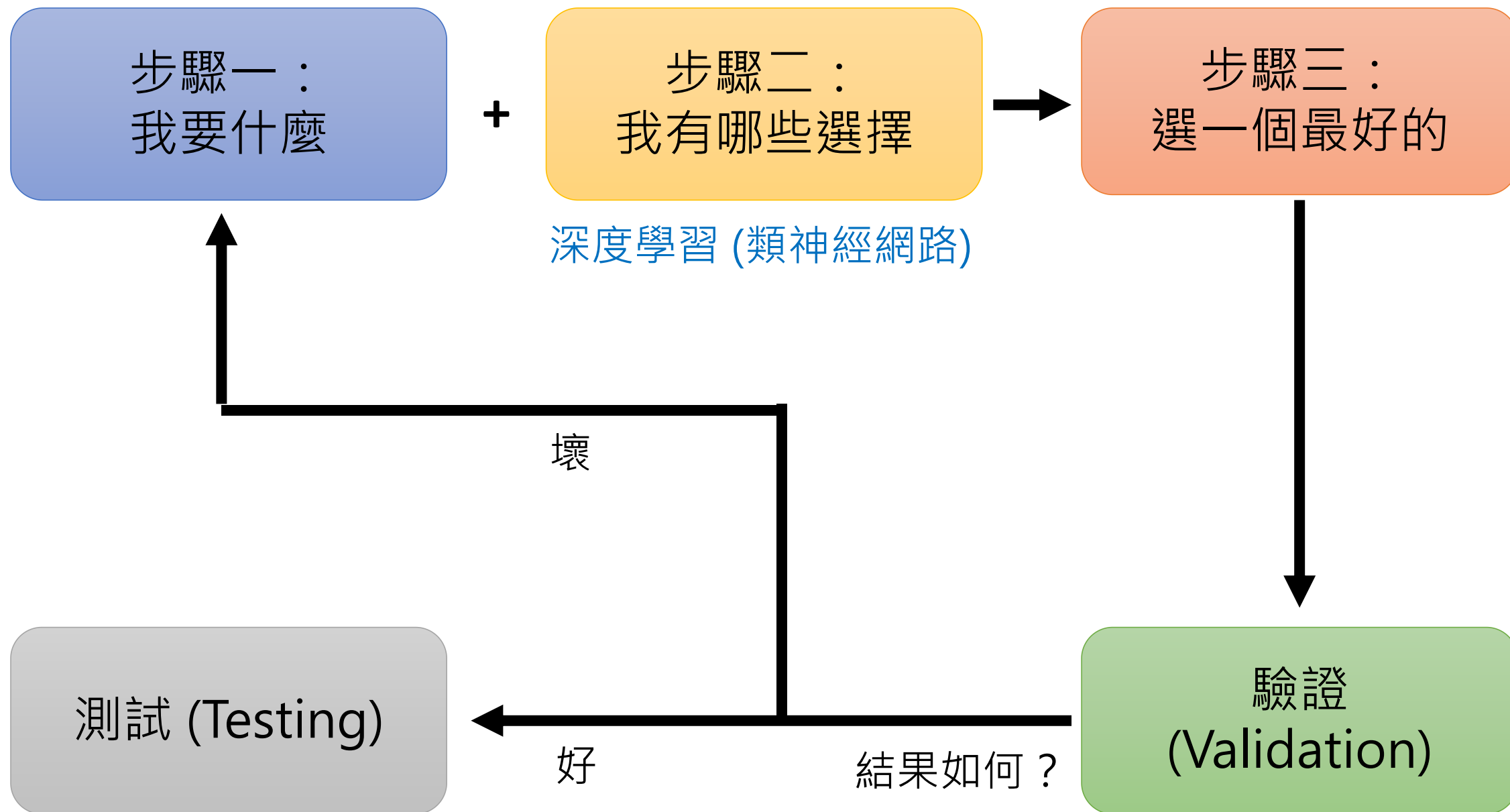


一堂課看懂 訓練類神經網路 的各種訣竅

李宏毅

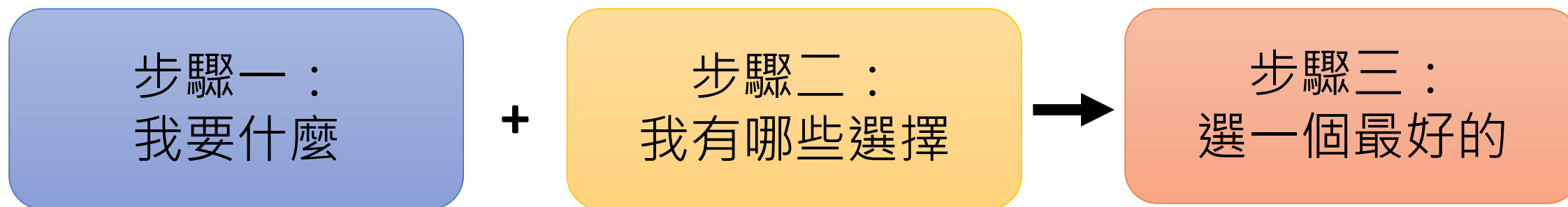




介紹各種訓練類神經網路常用技巧

- 聽到一個跟訓練類神經網路有關的方法時，你要這樣問自己

方法名	改了那一個步驟	帶來什麼好處
.....



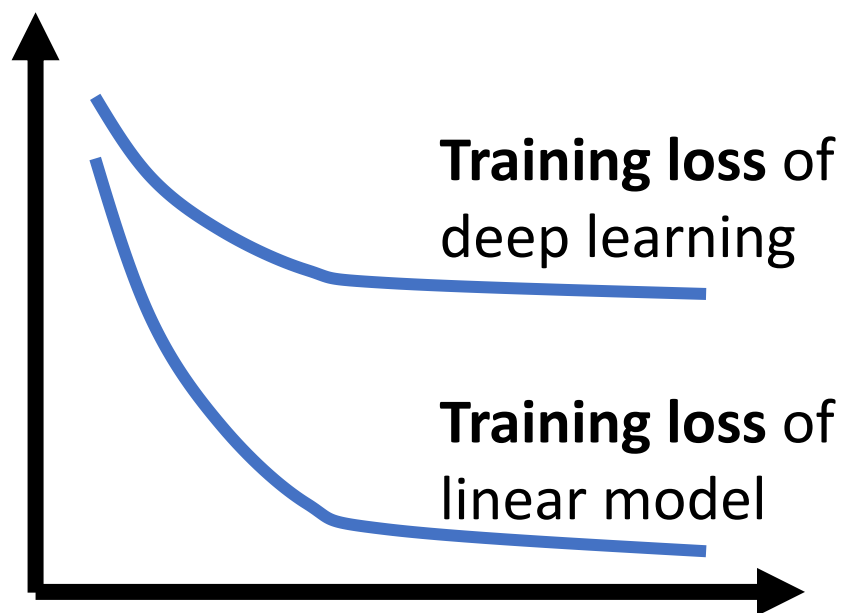
介紹各種訓練類神經網路常用技巧

- 聽到一個跟訓練類神經網路有關的方法時，你要這樣問自己

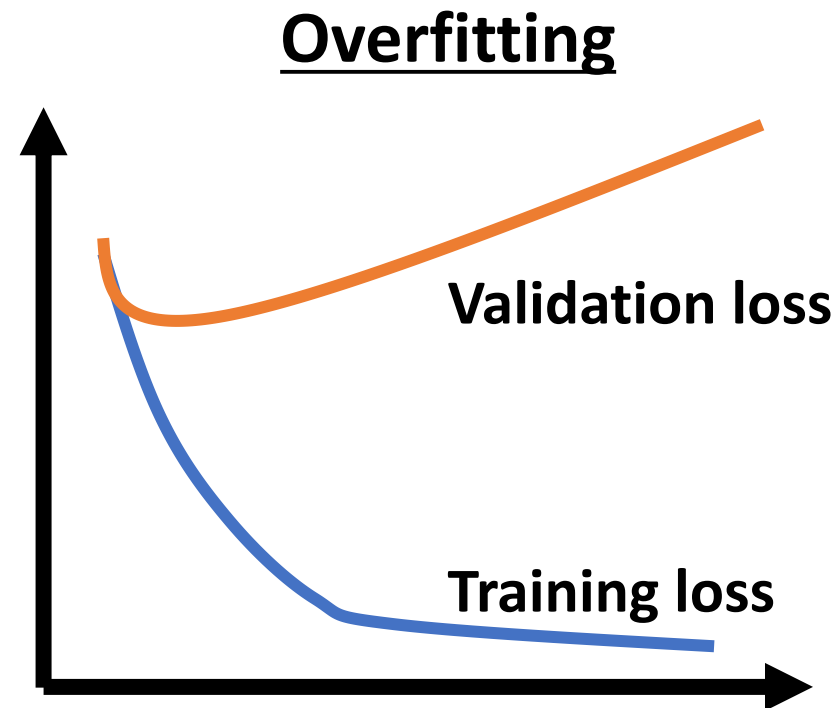
方法名	改了那一個步驟	帶來什麼好處
.....

- Better Optimization：更低的 Training Loss
- Better Generalization：更低的 Validation Loss
-

選擇合適的技巧

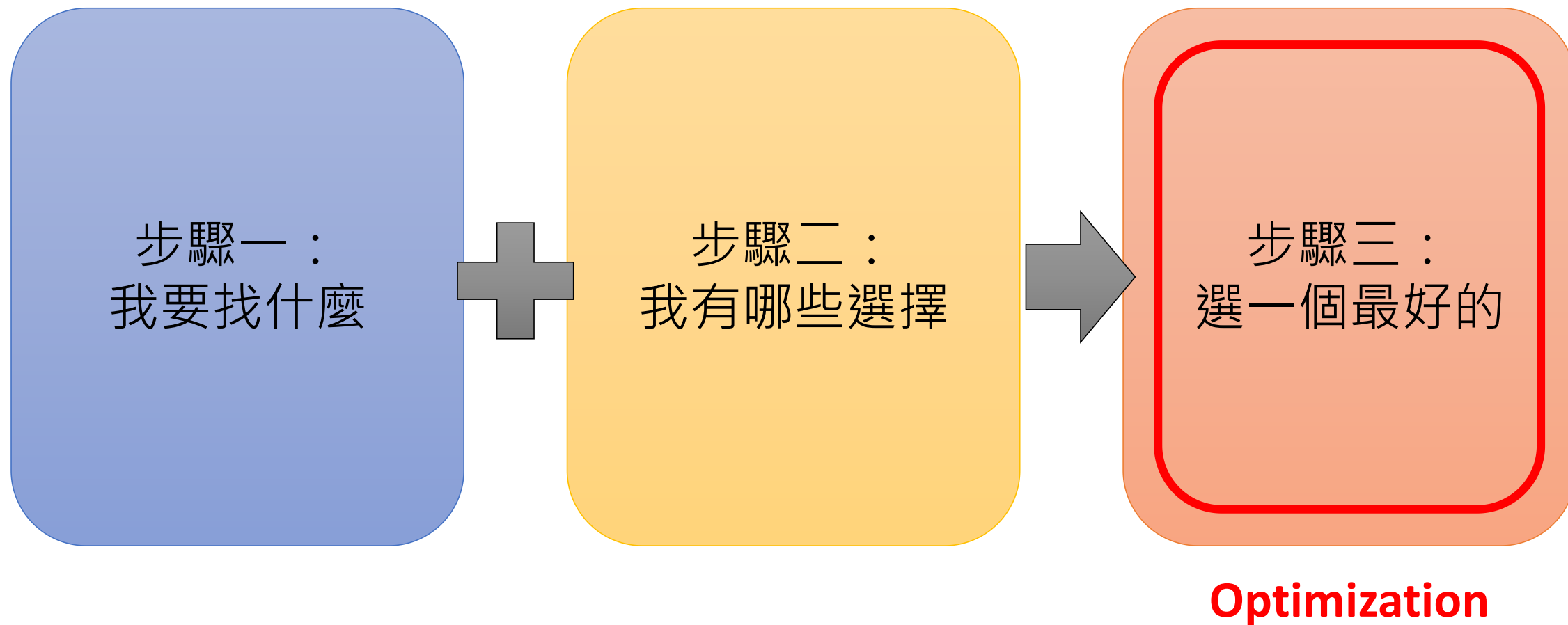


Better Optimization!



Better Generalization!

找函式步驟 3 + 1



Vanilla Gradient Descent

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} L(\boldsymbol{\theta})$$

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

➤ Compute gradient $\boldsymbol{g}^0 = \nabla L(\boldsymbol{\theta}^0)$

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

➤ Compute gradient $\boldsymbol{g}^1 = \nabla L(\boldsymbol{\theta}^1)$

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

➤ Compute gradient $\boldsymbol{g}^2 = \nabla L(\boldsymbol{\theta}^2)$

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

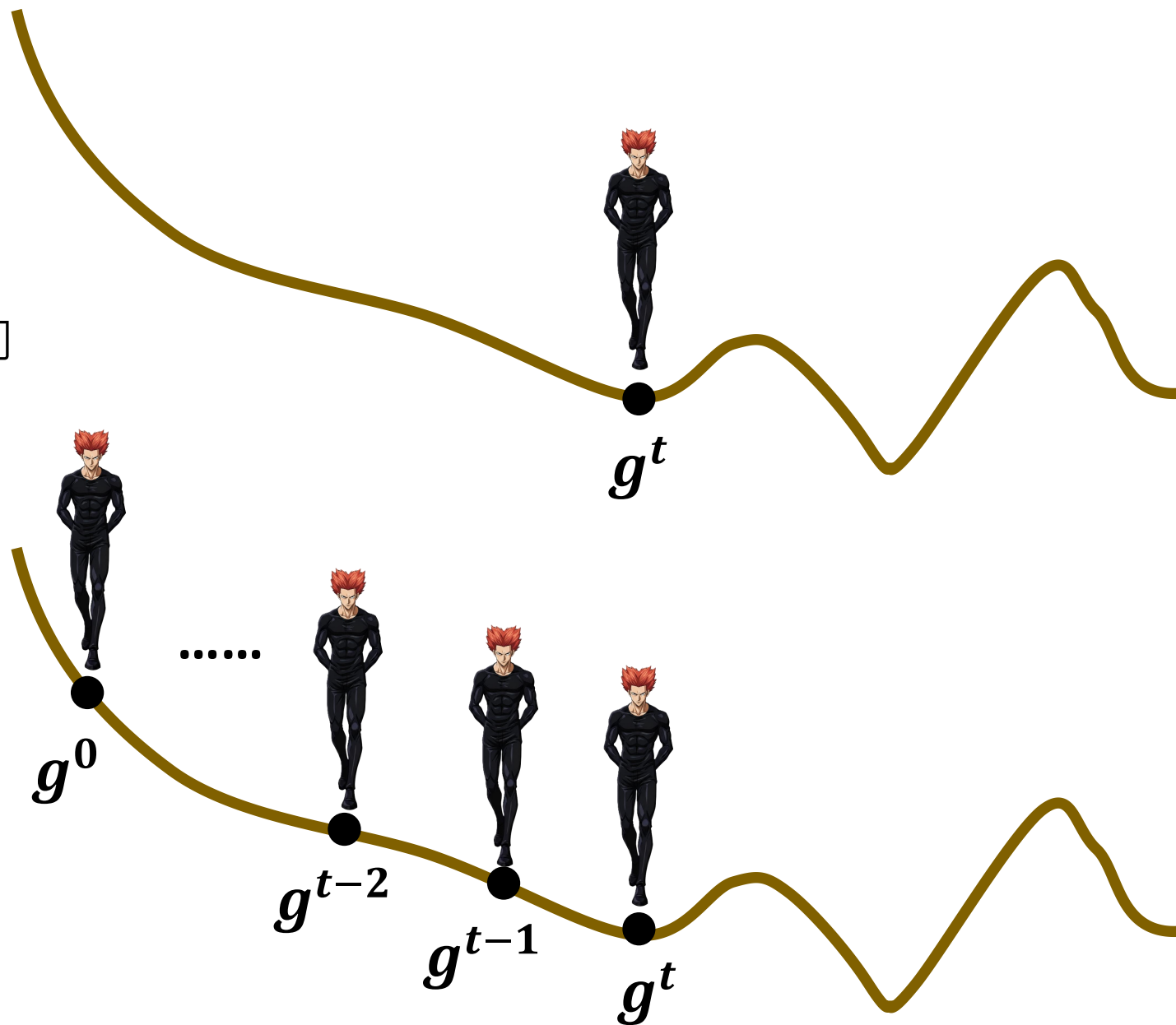
Optimizer

Vanilla Gradient Descent

根據當下算出來的 g^t 來決定方向

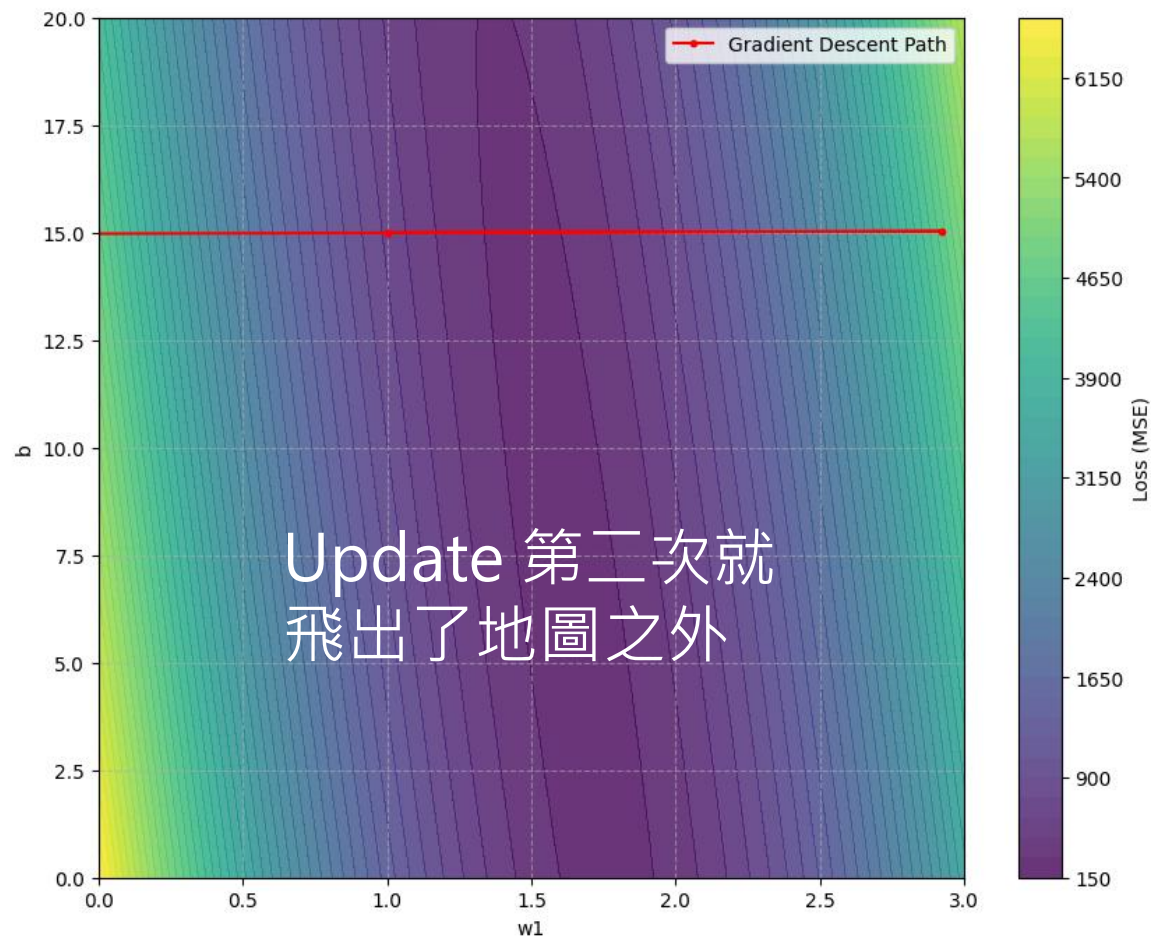
Gradient Descent + Optimizer

根據 $g^0, g^1, g^2, \dots, g^t$
一起來決定方向

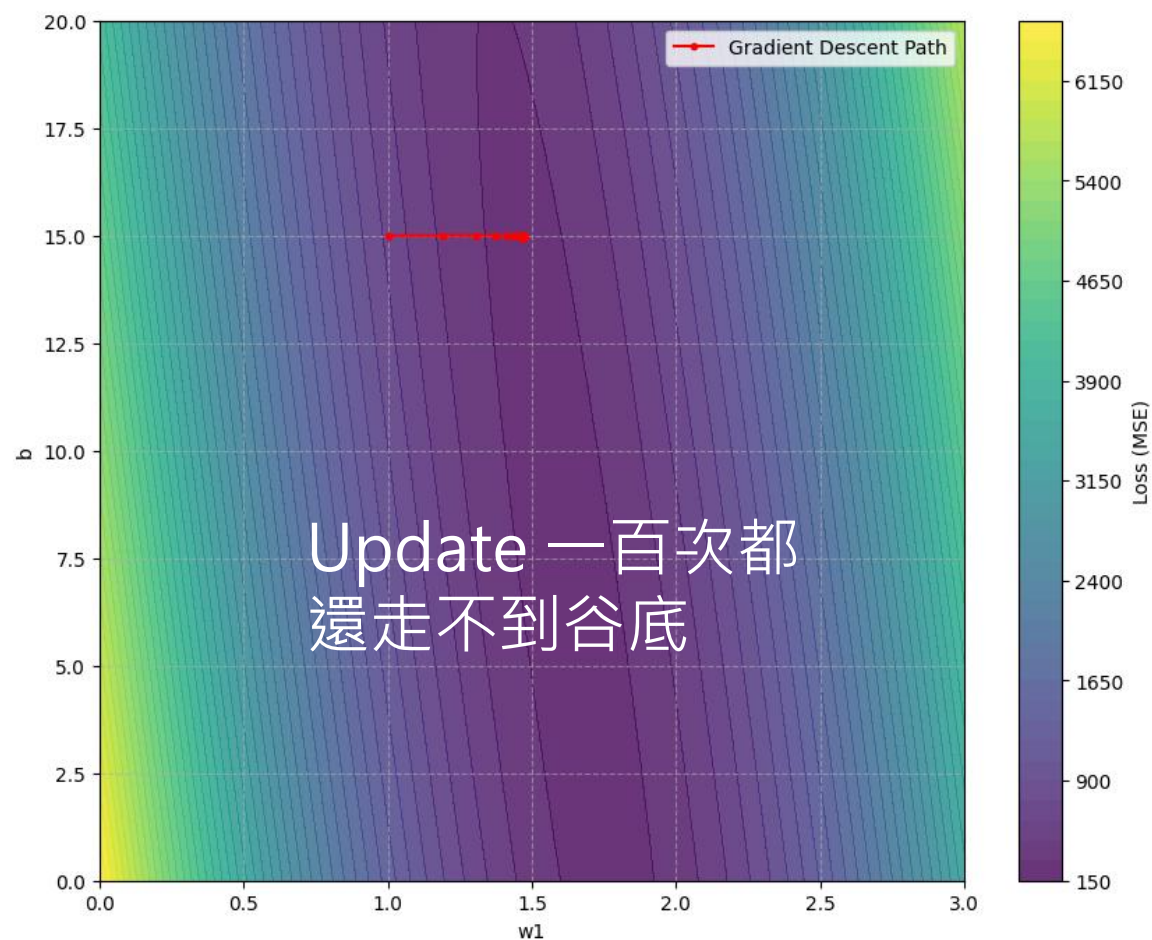


Learning Rate 實在是很難調 ...

$\eta = 0.001$

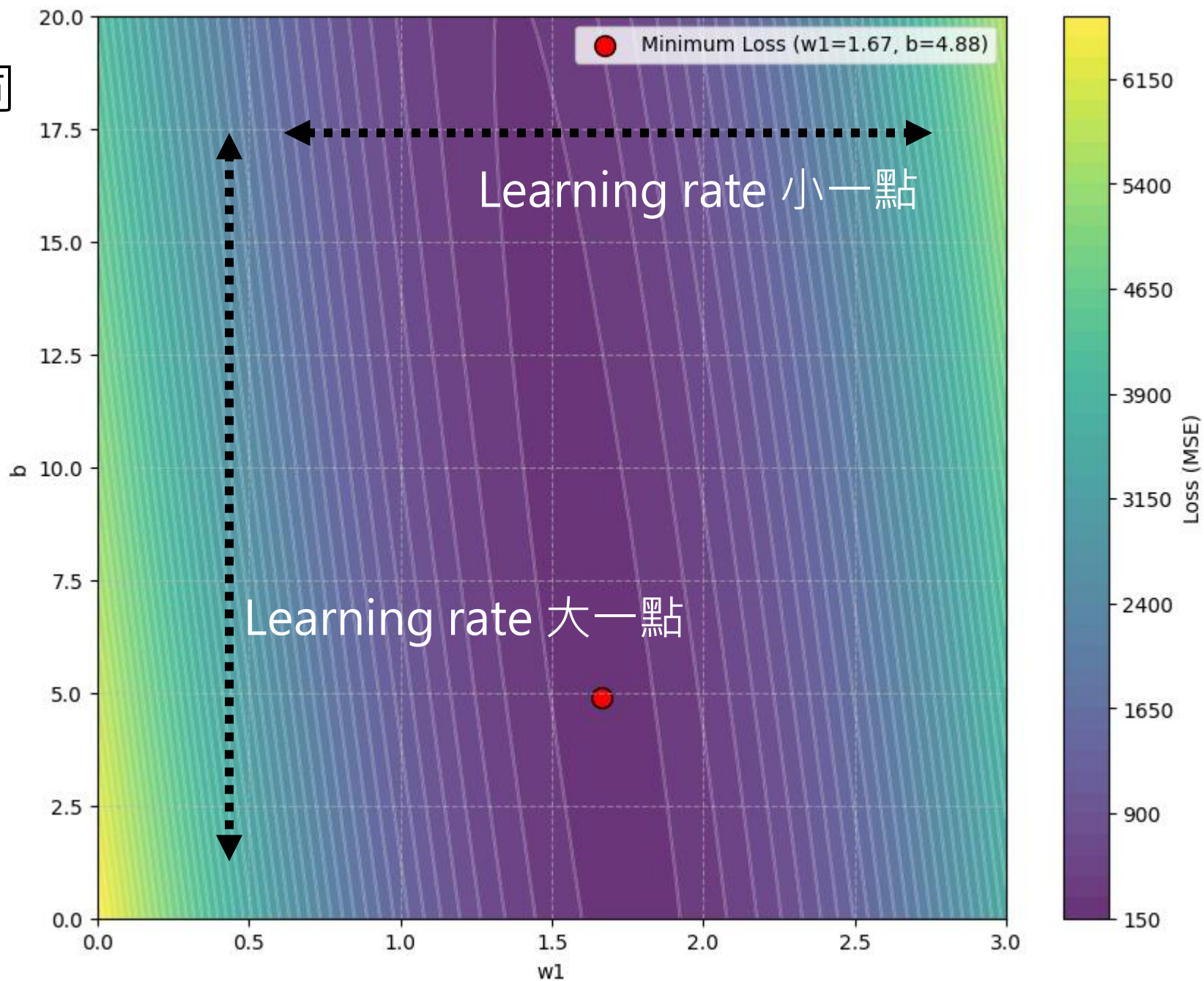


$\eta = 0.0001$

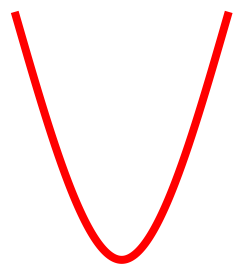


不同參數應該要有不同的 Learning Rate

怎麼知道那一個方向的 gradient 大、那一個小？



用過去的 Gradient 來決定 Learning Rate



Learning rate 小一點

$$g^0 = \begin{bmatrix} 500 \\ 0.4 \end{bmatrix}$$

$$g^1 = \begin{bmatrix} 432 \\ -0.3 \end{bmatrix}$$

$$g^2 = \begin{bmatrix} -211 \\ 0.2 \end{bmatrix}$$

$$g^3 = \begin{bmatrix} 139 \\ 0.1 \end{bmatrix}$$



Learning rate 大一點

Adagrad

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

For each dimension i : $\sigma_i^0 = \sqrt{(g_i^0)^2} = |g_i^0|$

$$\theta_i^1 \leftarrow \theta_i^0 - \frac{\eta}{\sigma_i^0} g_i^0$$

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

For each dimension i : $\sigma_i^1 = \sqrt{[(g_i^0)^2 + (g_i^1)^2]}$ Average?

$$\theta_i^2 \leftarrow \theta_i^1 - \frac{\eta}{\sigma_i^1} g_i^1$$

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

For each dimension i : $\sigma_i^2 = \sqrt{[(g_i^0)^2 + (g_i^1)^2 + (g_i^2)^2]}$

$$\theta_i^3 \leftarrow \theta_i^2 - \frac{\eta}{\sigma_i^2} g_i^2$$

⋮

- Compute gradient $\mathbf{g}^t = \nabla L(\boldsymbol{\theta}^t)$

For each dimension i : $\sigma_i^t = \sqrt{\sum_{i=0}^t (g_i^t)^2}$

$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t$$

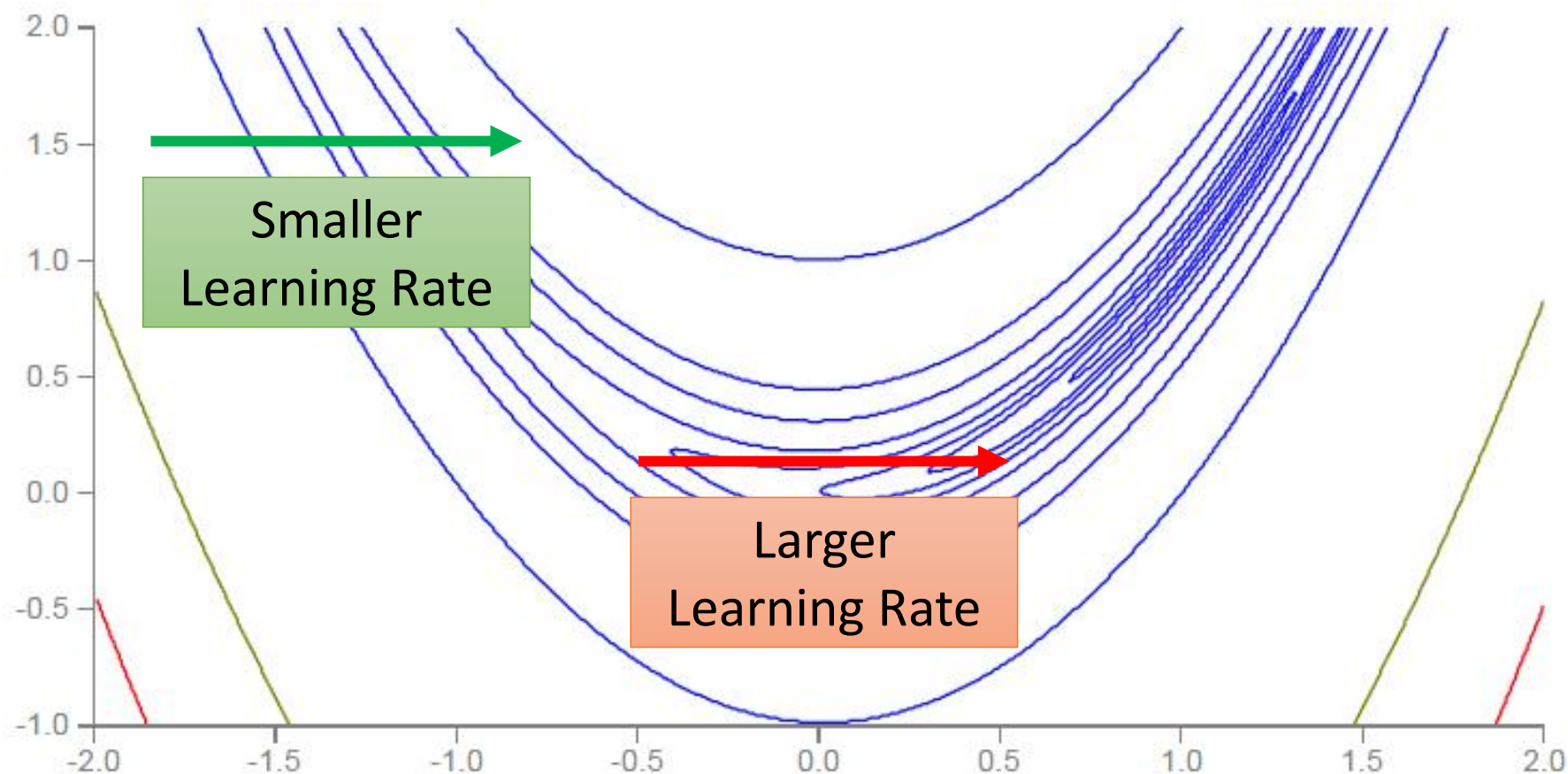
範例程式

連結：

<https://colab.research.google.com/drive/1XPIU-I77dXL9W74jnevmfb8K8XoEPRso?usp=sharing>



同一個參數的Gradient 大小不會一成不變



同一個參數的Gradient 大小不會一成不變

$$g^0 = \begin{bmatrix} 500 \\ 0.4 \end{bmatrix} \quad g^1 = \begin{bmatrix} 432 \\ -0.3 \end{bmatrix} \quad \dots \quad g^{t-1} = \begin{bmatrix} -0.1 \\ 229 \end{bmatrix} \quad g^t = \begin{bmatrix} 0.1 \\ 100 \end{bmatrix}$$

變大 變大



Adagrad: 全部平方加起來

RMSProp: 最近算出來的 gradient 給比較大的影響

RMSProp

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

For each dimension i : $\sigma_i^0 = \sqrt{(g_i^0)^2}$

$$\theta_i^1 \leftarrow \theta_i^0 - \frac{\eta}{\sigma_i^0} g_i^0$$

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

$$0 < \alpha < 1$$

For each dimension i : $\sigma_i^1 = \sqrt{\alpha(\sigma_i^0)^2 + (1 - \alpha)(g_i^1)^2}$

$$\theta_i^2 \leftarrow \theta_i^1 - \frac{\eta}{\sigma_i^1} g_i^1$$

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

For each dimension i : $\sigma_i^2 = \sqrt{\alpha(\sigma_i^1)^2 + (1 - \alpha)(g_i^2)^2}$

$$\theta_i^3 \leftarrow \theta_i^2 - \frac{\eta}{\sigma_i^2} g_i^2$$

⋮

- Compute gradient $\mathbf{g}^t = \nabla L(\boldsymbol{\theta}^t)$

For each dimension i : $\sigma_i^t = \sqrt{\alpha(\sigma_i^{t-1})^2 + (1 - \alpha)(g_i^t)^2}$

$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t$$

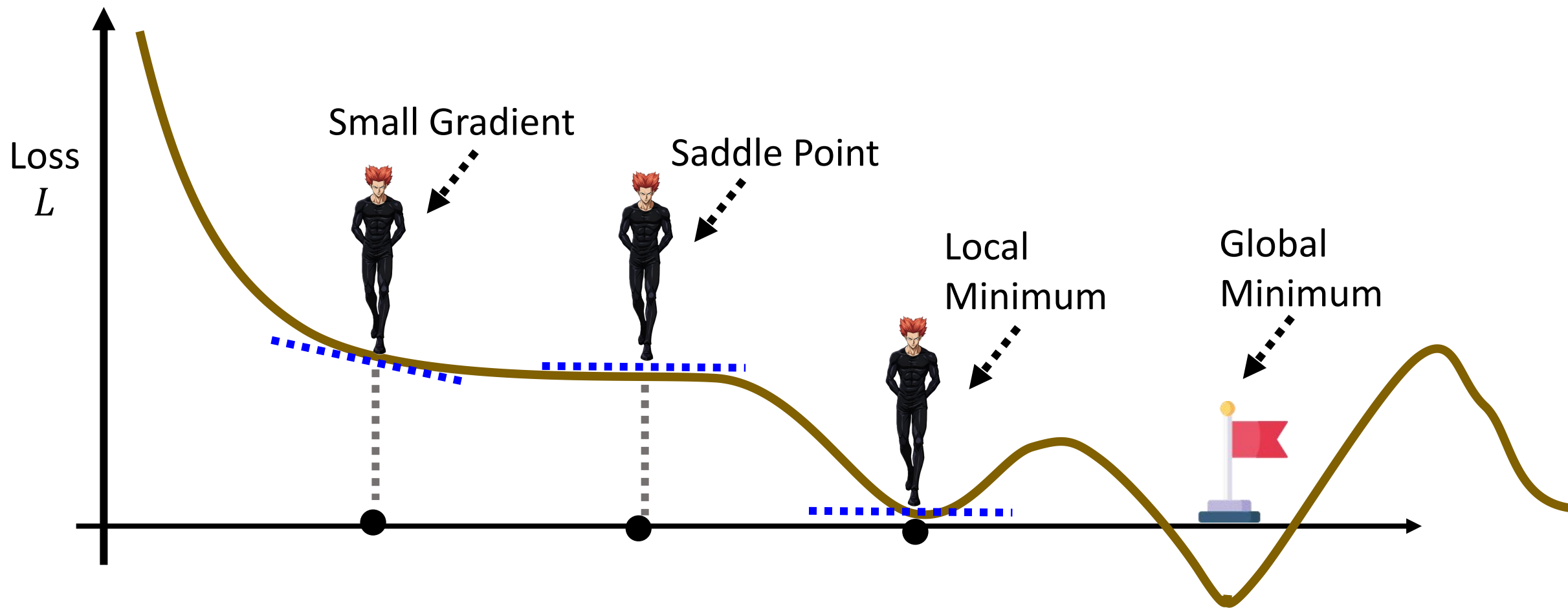
範例程式

連結：

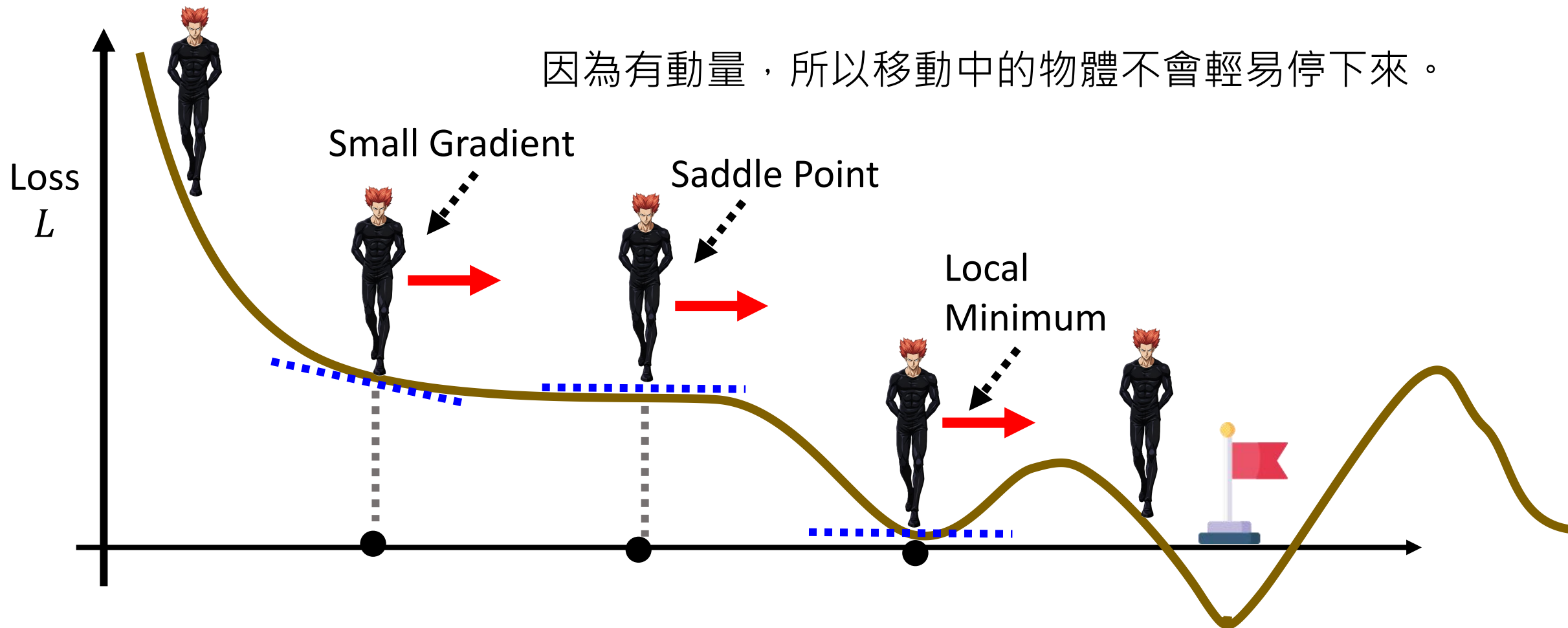
<https://colab.research.google.com/drive/1XPIU-I77dXL9W74jnevmfb8K8XoEPRso?usp=sharing>



Optimization 會在 Gradient 很小時停止



考慮動量 (Momentum)



Momentum

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

For each dimension i : $m_i^0 = g_i^0$

$$\theta_i^1 \leftarrow \theta_i^0 - \eta m_i^0$$

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

For each dimension i : $m_i^1 = g_i^0 + g_i^1$

$$\theta_i^2 \leftarrow \theta_i^1 - \eta m_i^1$$

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

For each dimension i : $m_i^2 = g_i^0 + g_i^1 + g_i^2$

$$\theta_i^3 \leftarrow \theta_i^2 - \eta m_i^2$$

⋮

- Compute gradient $\mathbf{g}^t = \nabla L(\boldsymbol{\theta}^t)$

For each dimension i : $m_i^t = g_i^0 + g_i^1 + g_i^2 + \dots + g_i^t$

$$\theta_i^{t+1} \leftarrow \theta_i^t - \eta m_i^t$$

Momentum (這不是經典的 Momentum)

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

For each dimension i : $m_i^0 = g_i^0$

$$\theta_i^1 \leftarrow \theta_i^0 - \eta m_i^0$$

➤ Compute gradient $\mathbf{g}^1 = \nabla L(\boldsymbol{\theta}^1)$ $0 < \beta < 1$

For each dimension i : $m_i^1 = \beta m_i^0 + (1 - \beta) g_i^1$

$$\theta_i^2 \leftarrow \theta_i^1 - \eta m_i^1$$

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

For each dimension i : $m_i^2 = \beta m_i^1 + (1 - \beta) g_i^2$

$$\theta_i^3 \leftarrow \theta_i^2 - \eta m_i^2$$

⋮

➤ Compute gradient $\mathbf{g}^t = \nabla L(\boldsymbol{\theta}^t)$

For each dimension i : $m_i^t = \beta m_i^{t-1} + (1 - \beta) g_i^t$

$$\theta_i^{t+1} \leftarrow \theta_i^t - \eta m_i^t$$

範例程式

連結：

<https://colab.research.google.com/drive/1XPIU-I77dXL9W74jnevmfb8K8XoEPRso?usp=sharing>



Adam: RMSProp + Momentum

- Compute gradient $\mathbf{g}^t = \nabla L(\boldsymbol{\theta}^t)$

Momentum

For each dimension i : $\mathbf{m}_i^t = \beta \mathbf{m}_i^{t-1} + (1 - \beta) \mathbf{g}_i^t$ $\theta_i^{t+1} \leftarrow \theta_i^t - \eta \mathbf{m}_i^t$

RMSProp

For each dimension i : $\sigma_i^t = \sqrt{\alpha (\sigma_i^{t-1})^2 + (1 - \alpha) (\mathbf{g}_i^t)^2}$ $\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} \mathbf{g}_i^t$

Adam

$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} \mathbf{m}_i^t$$

Adam has bias-corrected terms,
which are omitted here for simplicity.