

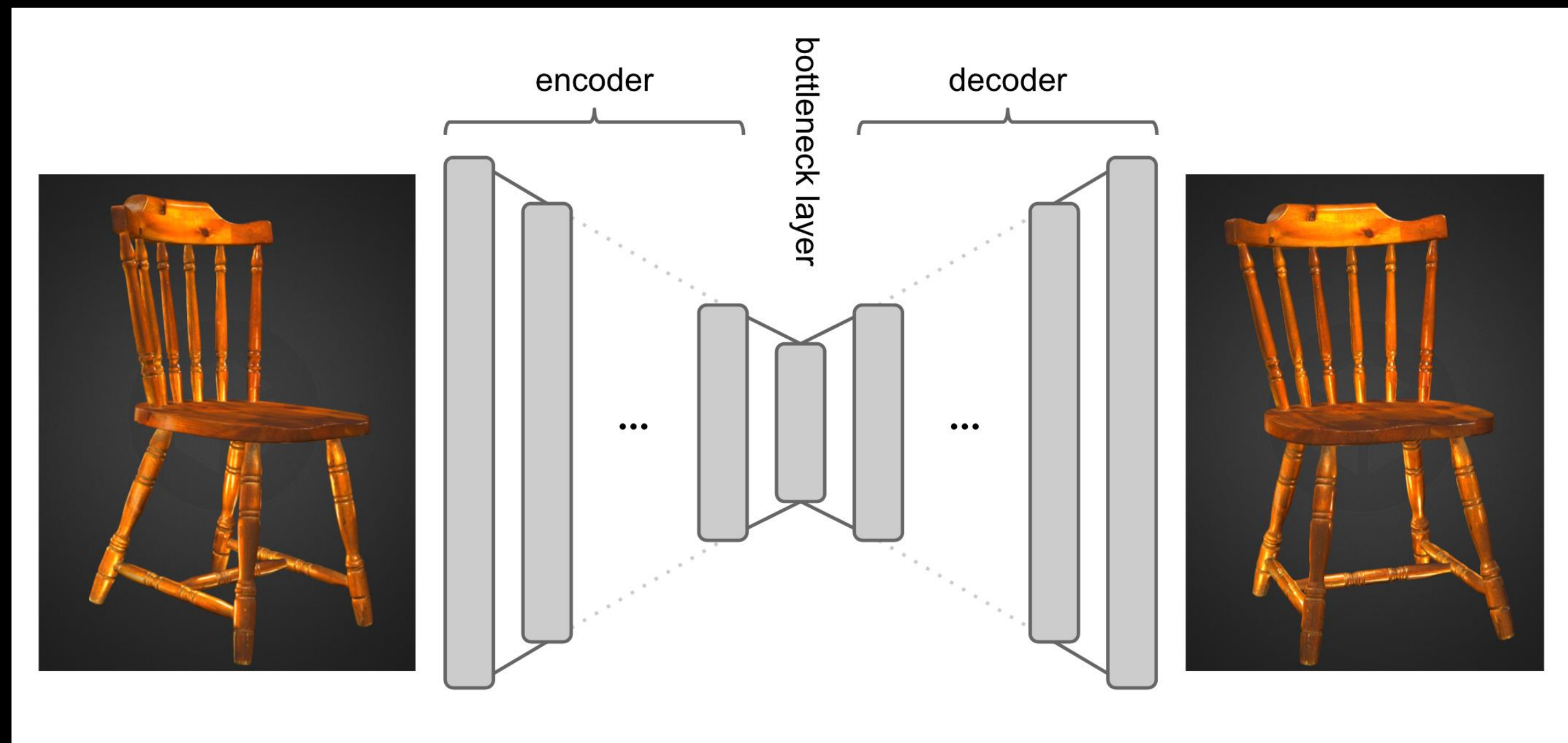
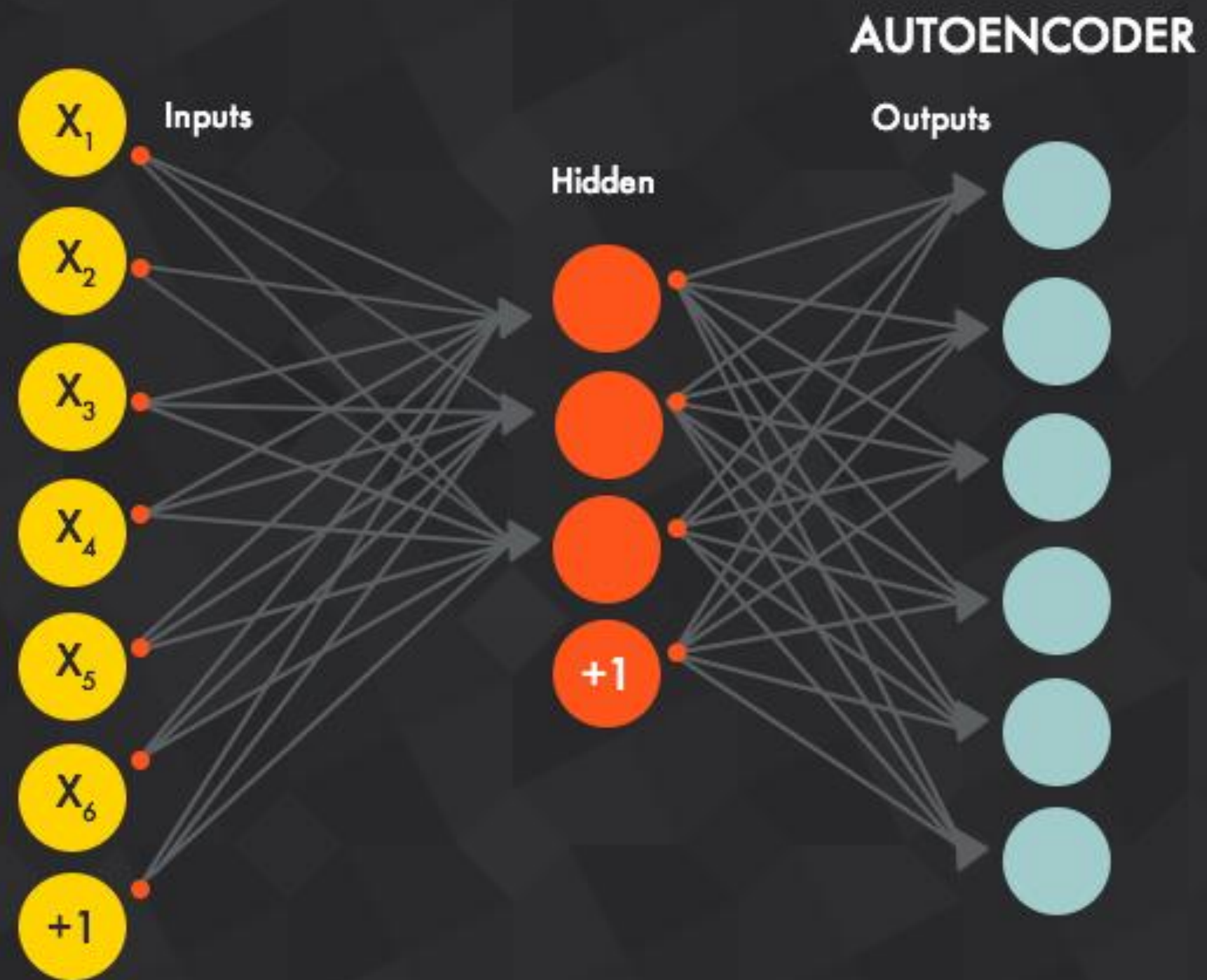
A blurred city street scene with a pink and blue banner overlay. The background shows a busy street with pedestrians and buildings. The banner is a horizontal strip with a pink main body and blue triangular ends on the left and right. The text "朝向最佳解邁進!" is written in white on the pink part of the banner.

朝向最佳解邁進!

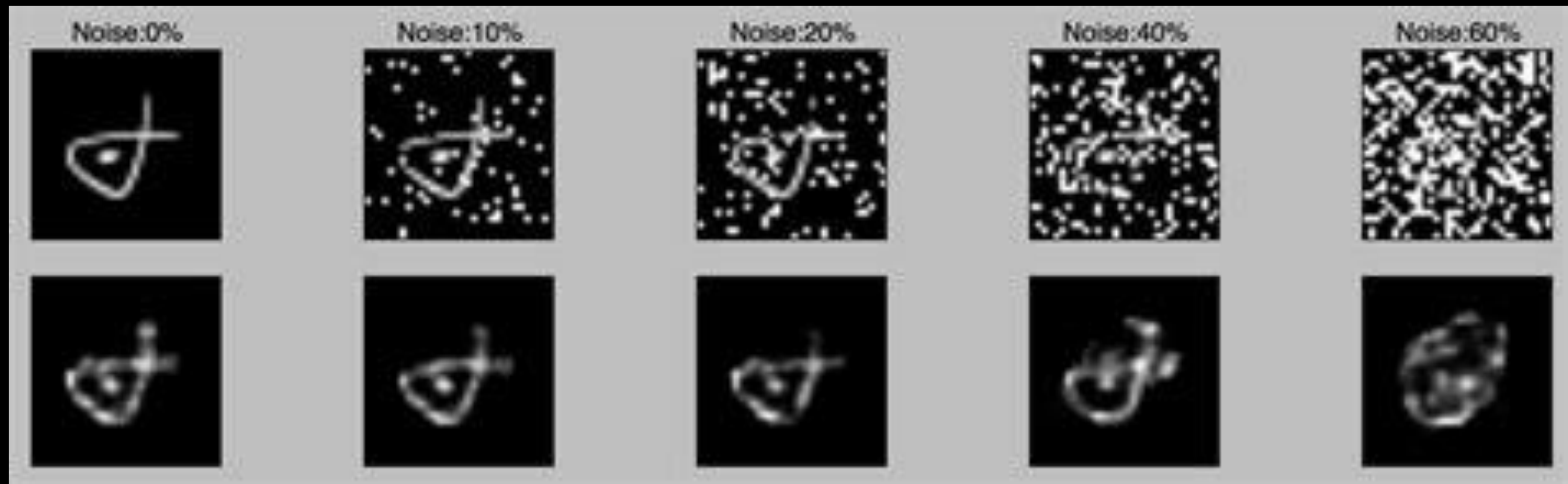
自動編碼器

Auto-encoder

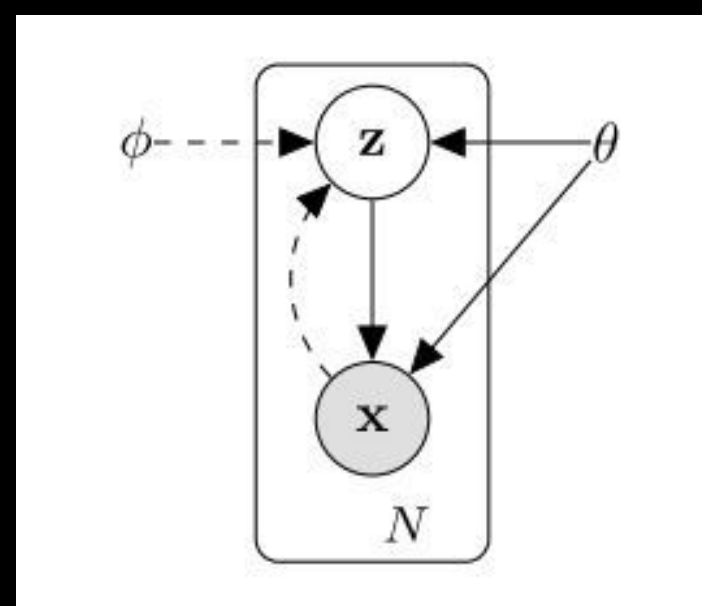
透過壓縮來尋找特徵



Denoising autoencoder

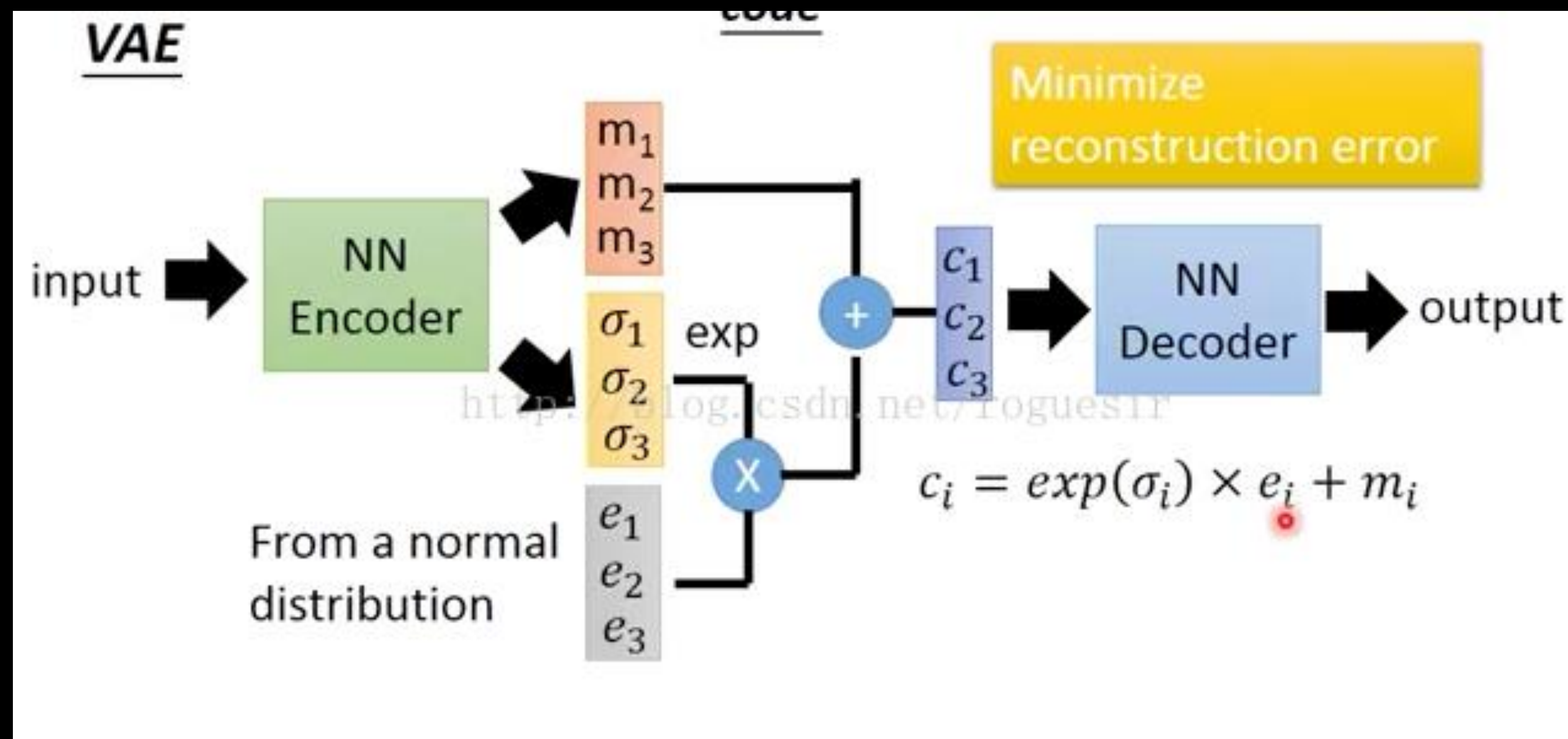


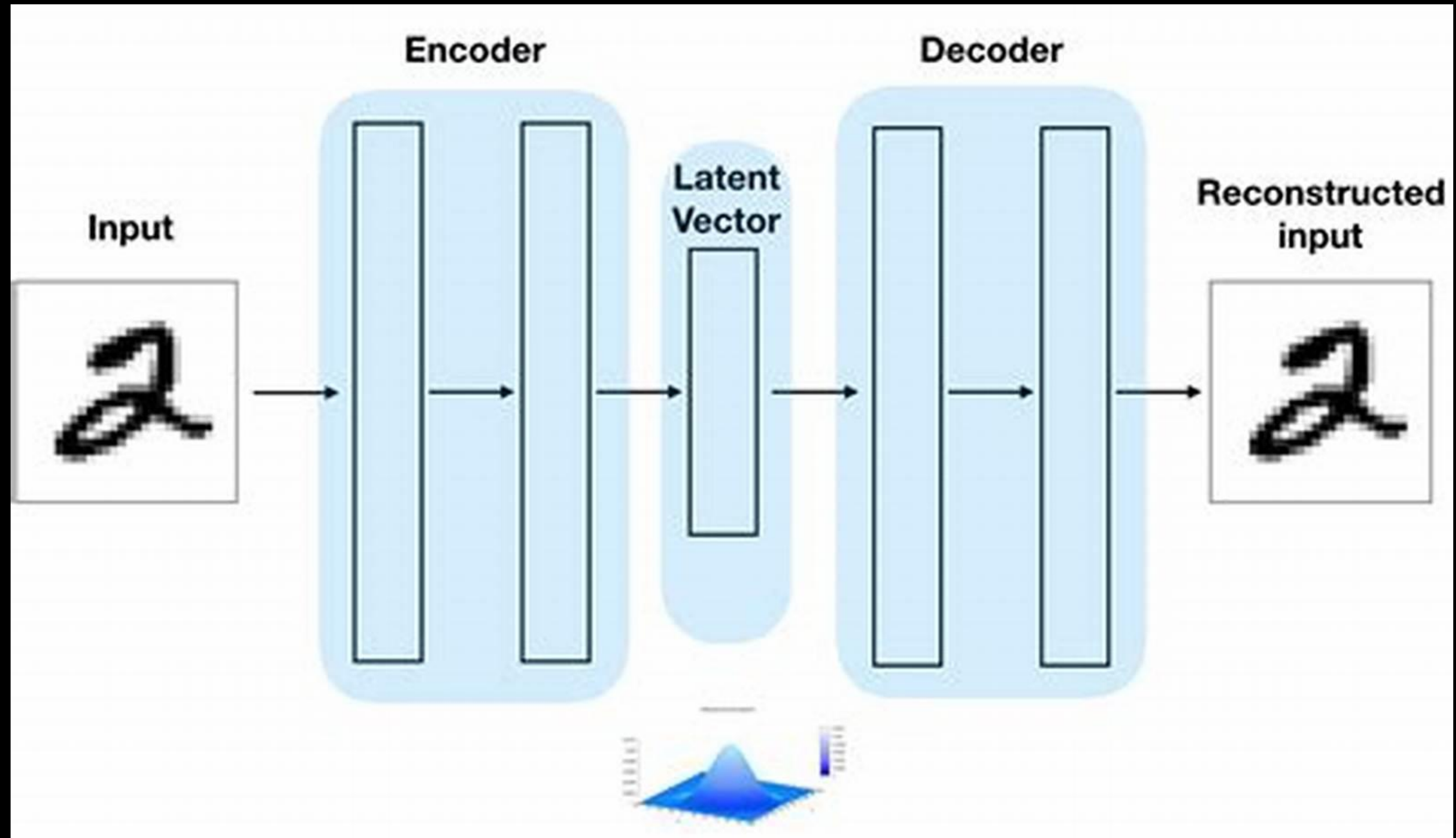
變分自編碼器 (VAE)








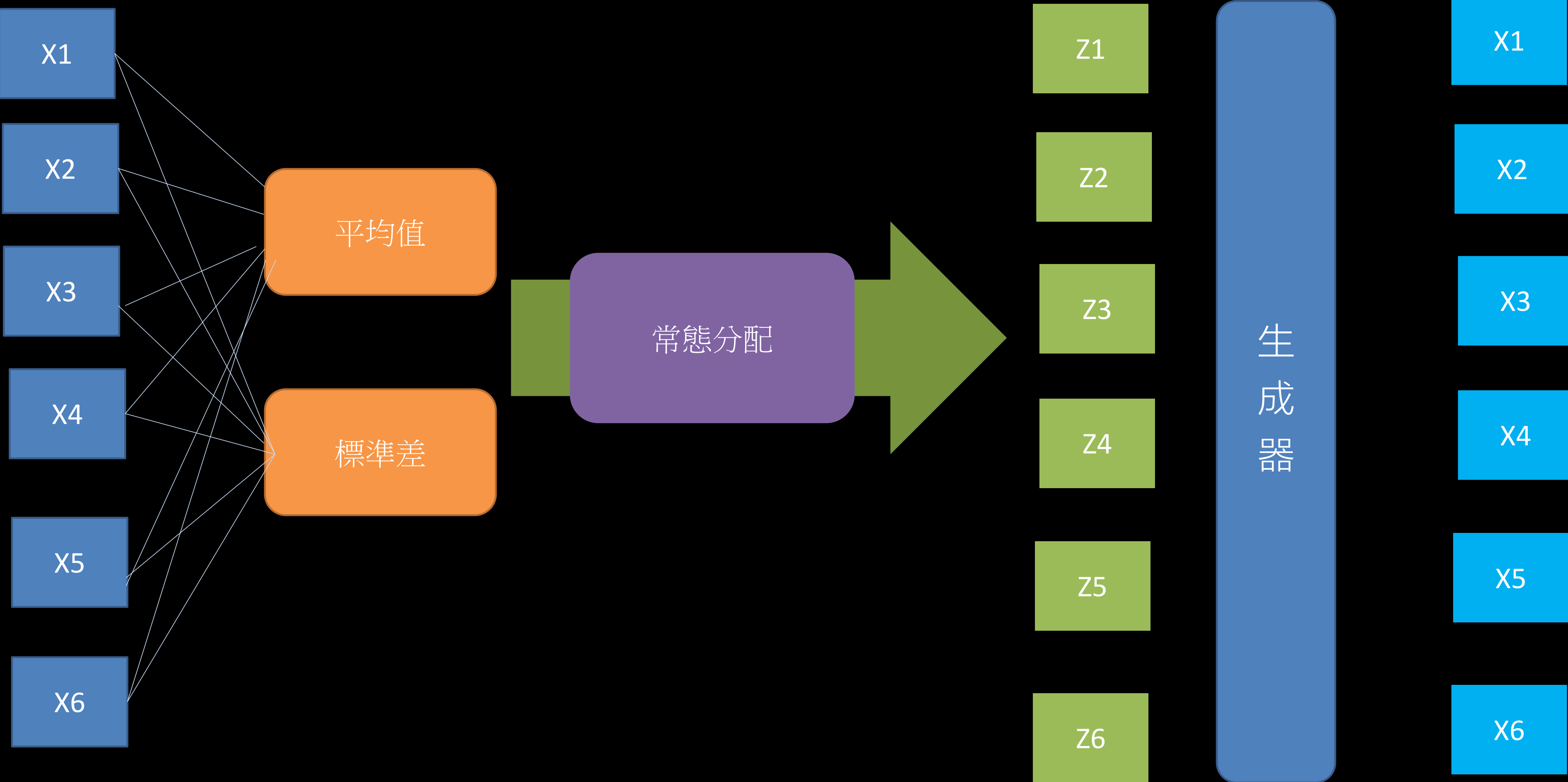
觀測到的資料是 x ，而由隱變數 z 產生，由 $z \rightarrow x$ 是生成模型，從自編碼器 (auto-encoder) 的角度來看，就是解碼器；而 $x \rightarrow z$ 由是識別模型 (recognition model)，類似於自編碼器的編碼器。

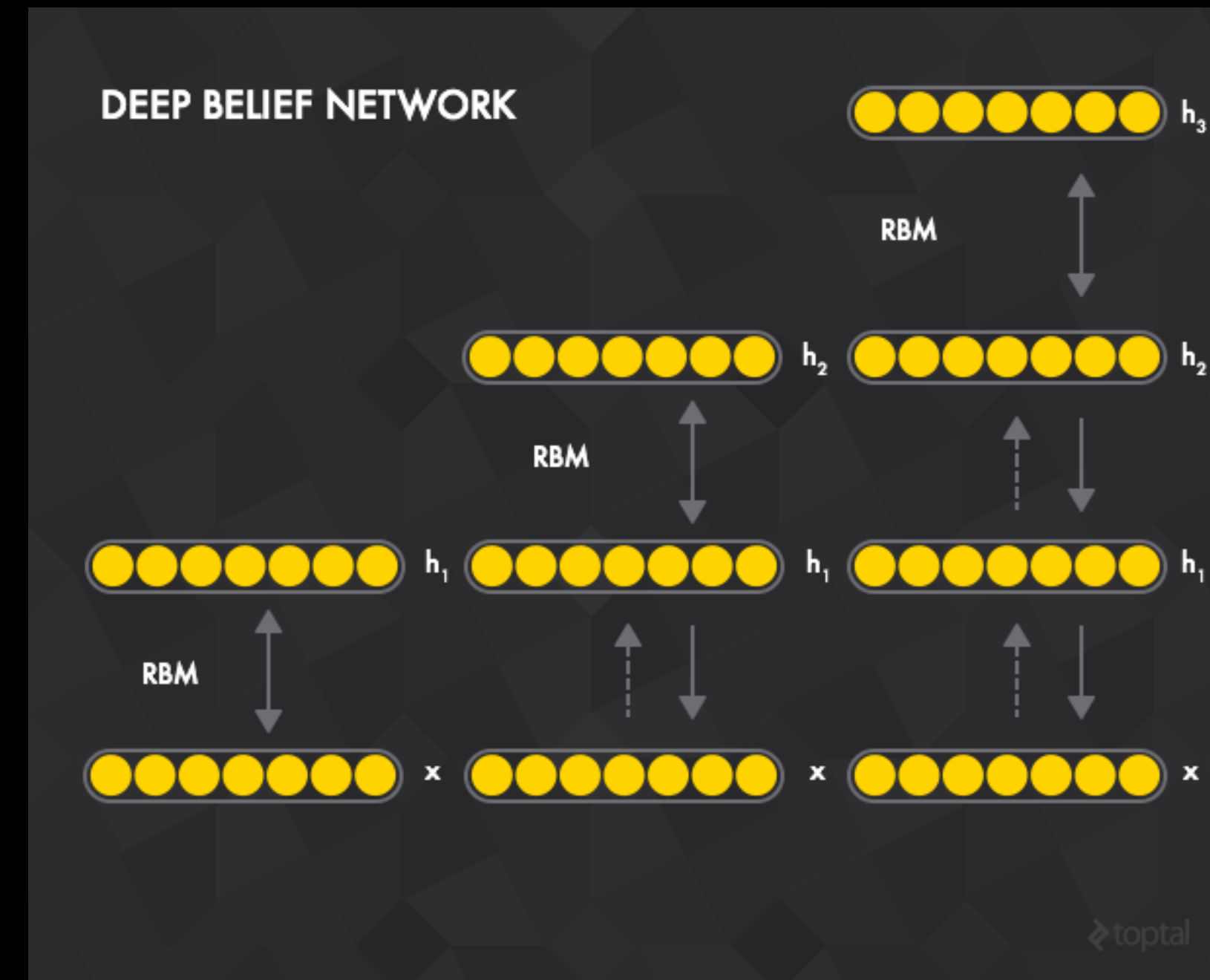
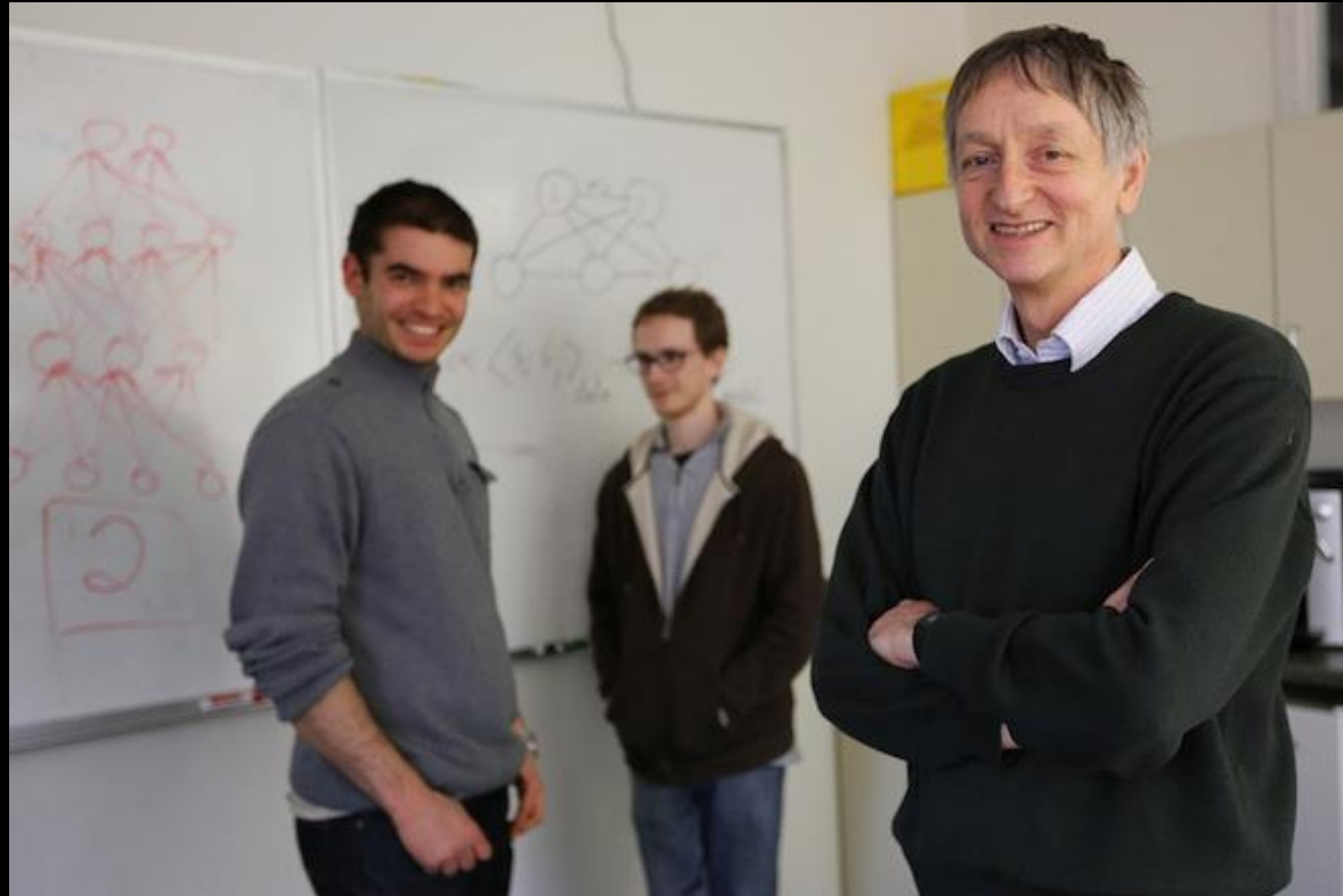
學習圖像的密度函數 (PDF)
兩個分佈的相似程度，一般採用KL散度





Input image	2-D latent space	5-D latent space	10-D latent space	20-D latent space
				



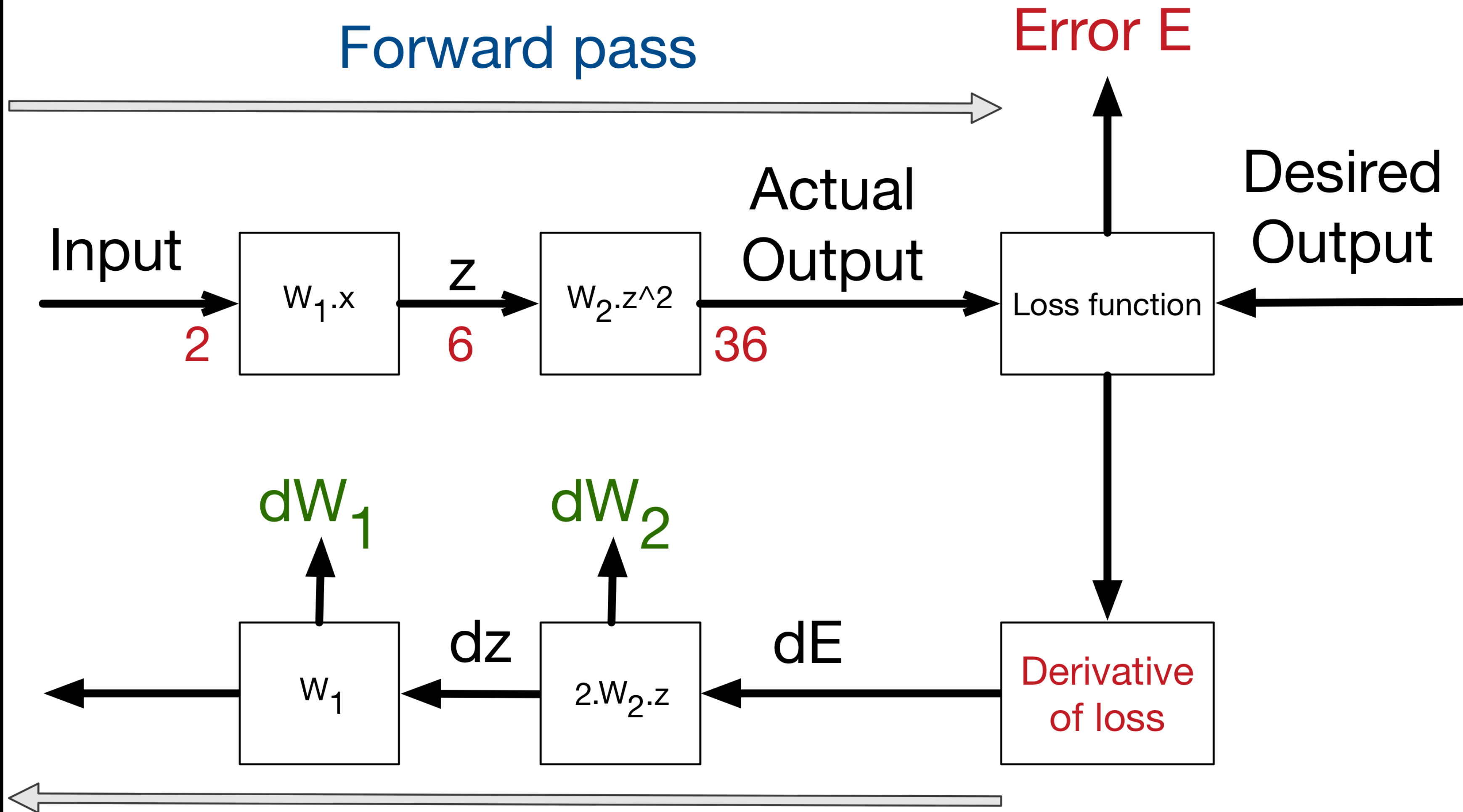


A fast learning algorithm for deep belief nets: 2006

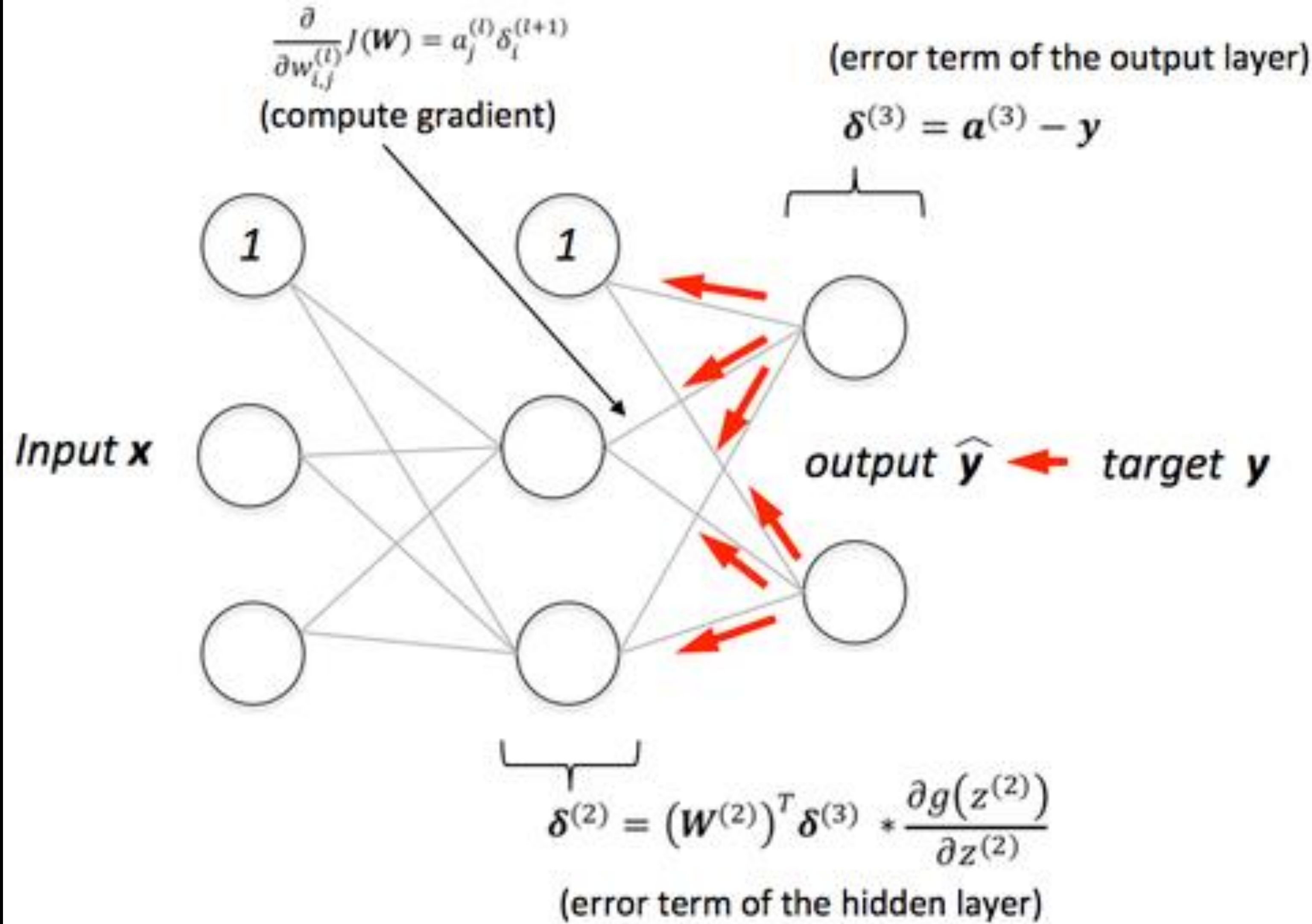
To RecogniZe Shapes, First Learn to Generate Images: 2006

Reducing the dimensionality of data with neural networks : 2006

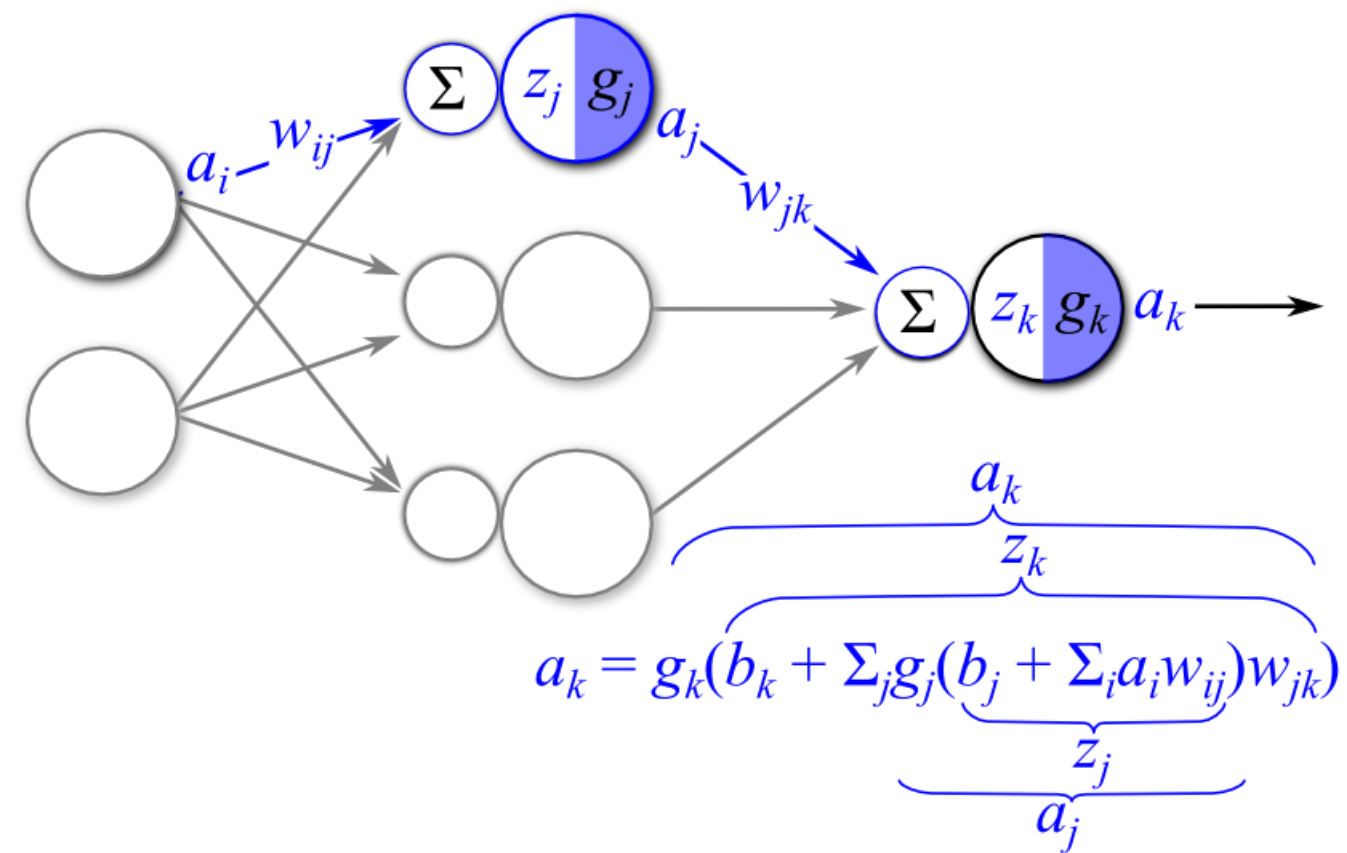
Forward pass



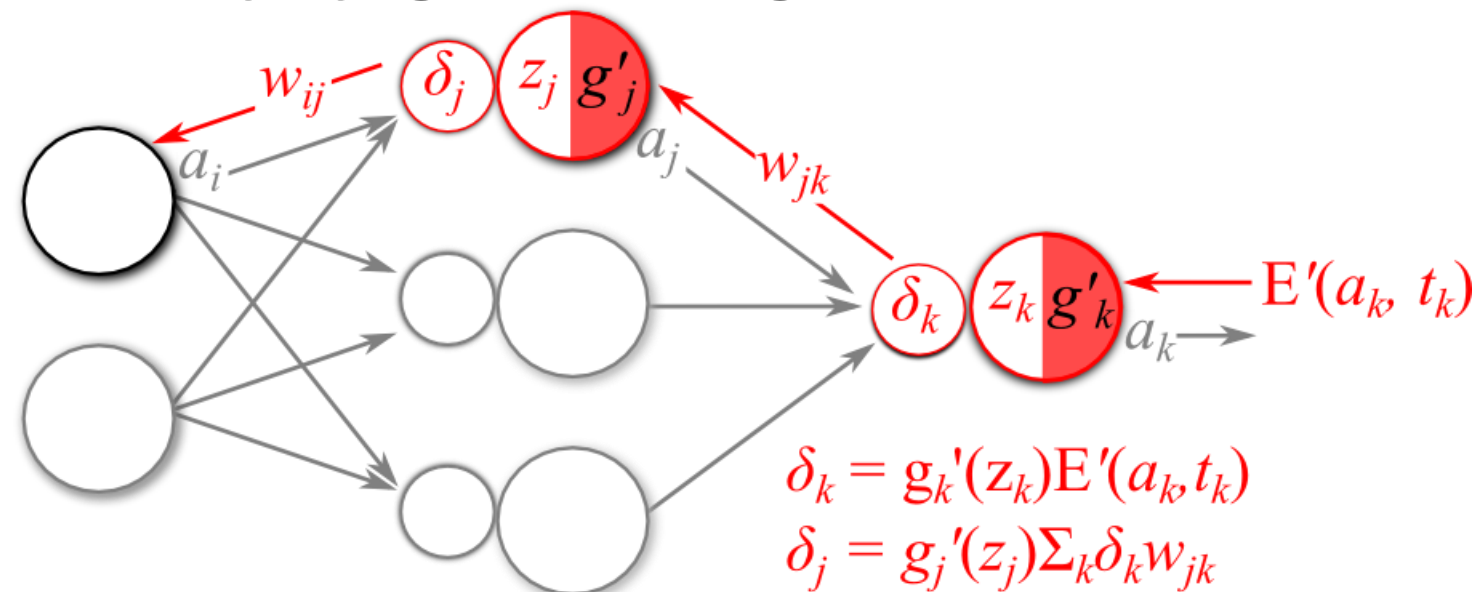
Back-propagate error



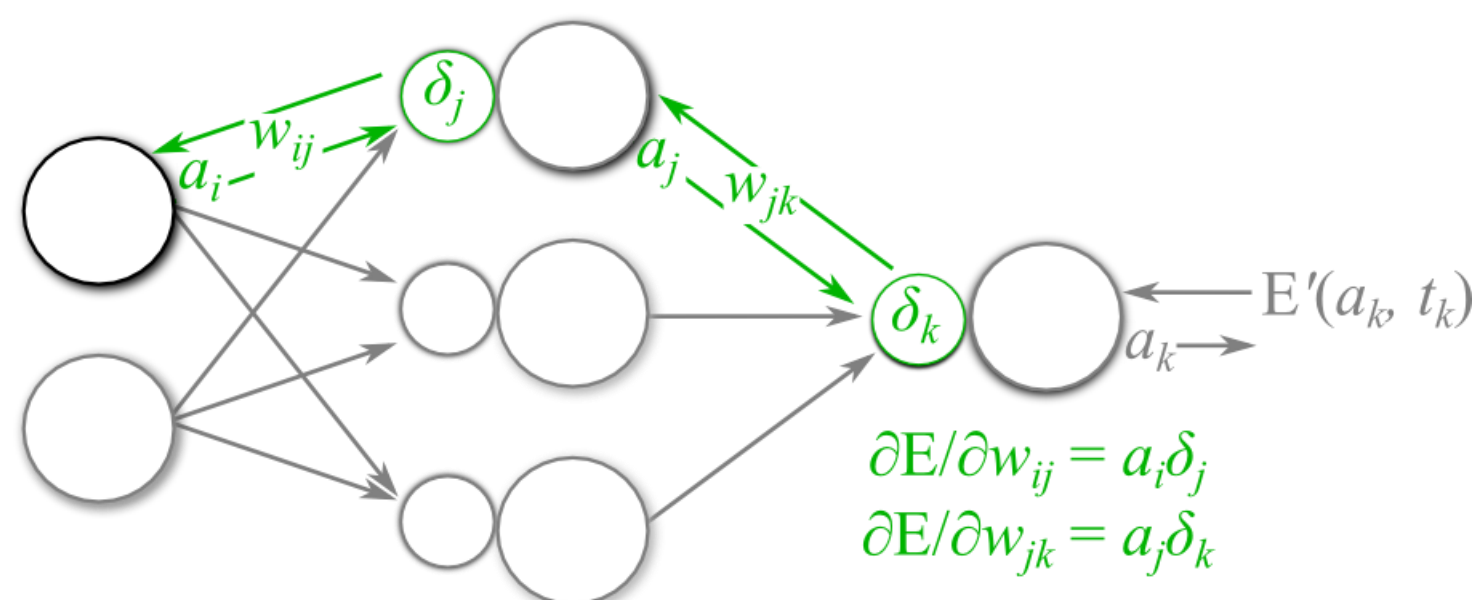
I. Forward-propagate Input Signal



II. Back-propagate Error Signals



III. Calculate Parameter Gradients



IV. Update Parameters

$$w_{ij} = w_{ij} - \eta (\frac{\partial E}{\partial w_{ij}})$$

$$w_{jk} = w_{jk} - \eta (\frac{\partial E}{\partial w_{jk}})$$

for learning rate η

順向傳導

逆向傳導

計算圖

Scalar	Vector	Matrix	Tensor
1	$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$	$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$	$\begin{bmatrix} \begin{bmatrix} 1 & 2 \end{bmatrix} & \begin{bmatrix} 3 & 2 \end{bmatrix} \\ \begin{bmatrix} 1 & 7 \end{bmatrix} & \begin{bmatrix} 5 & 4 \end{bmatrix} \end{bmatrix}$

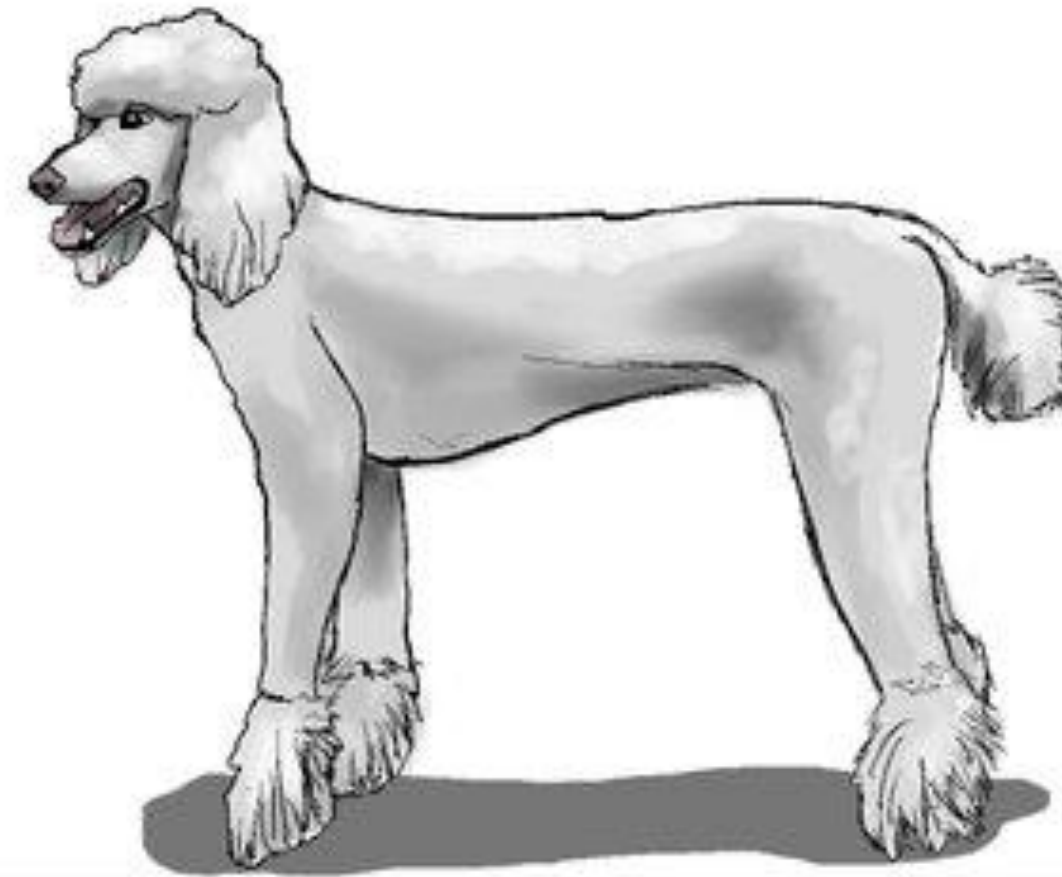
Scalar



Vector



Matrix



Tensor



給定待最佳化的模型參數 $\theta \in \mathbb{R}^d$

損失函數 $J(\theta)$

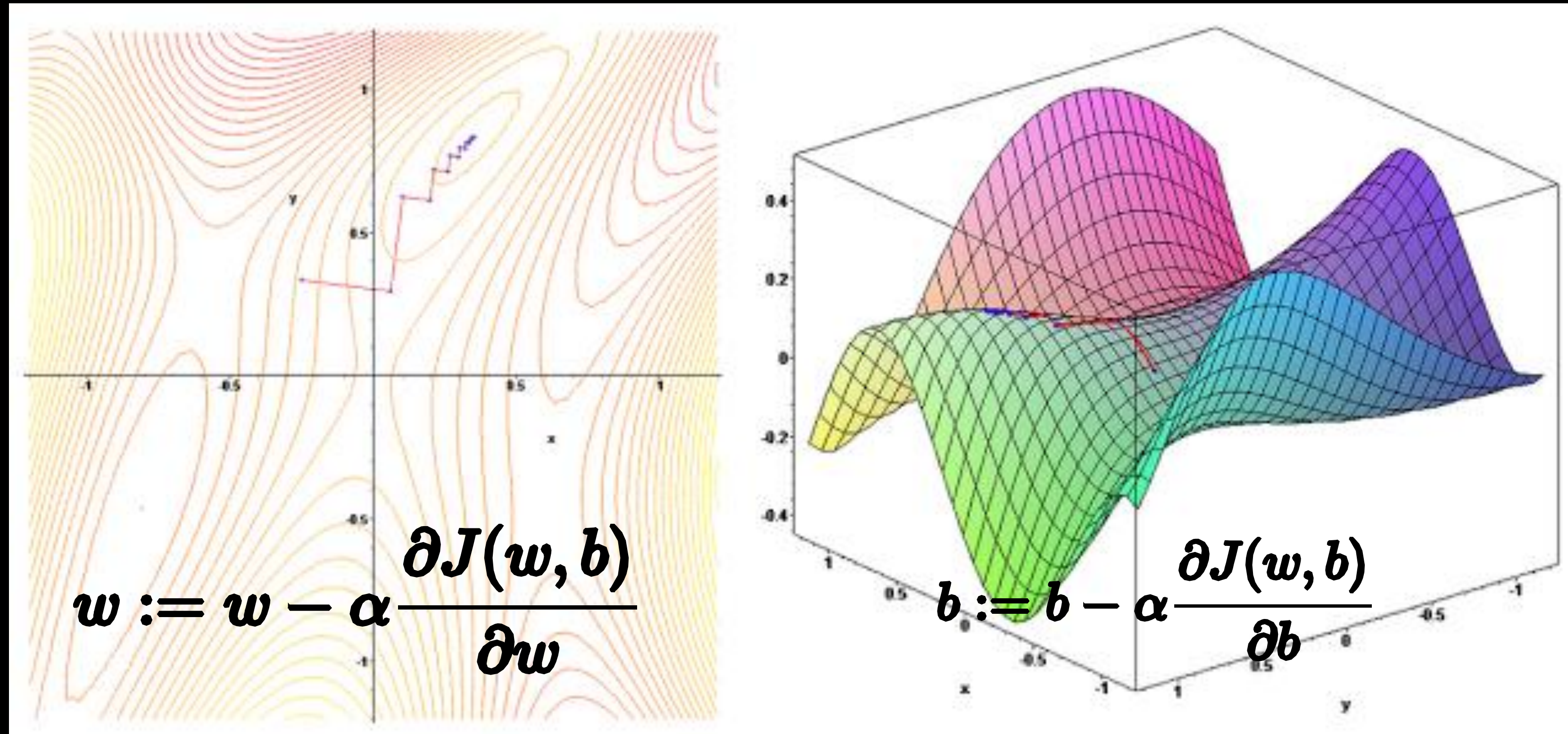
沿著 $\nabla_{\theta} J(\theta)$ 梯度向下的方向來更新 θ

學習速率LR決定了每一時刻的更新步長

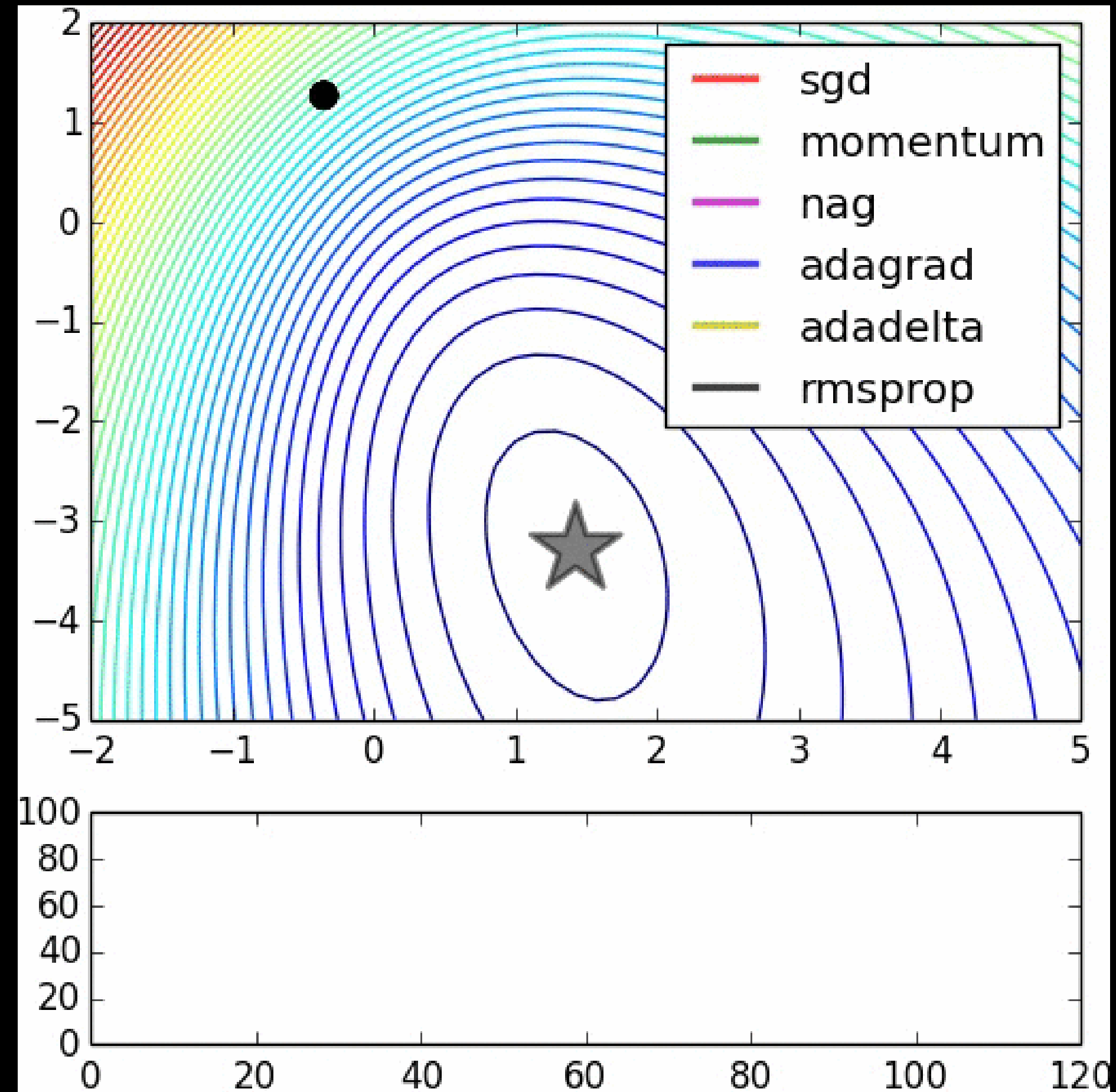
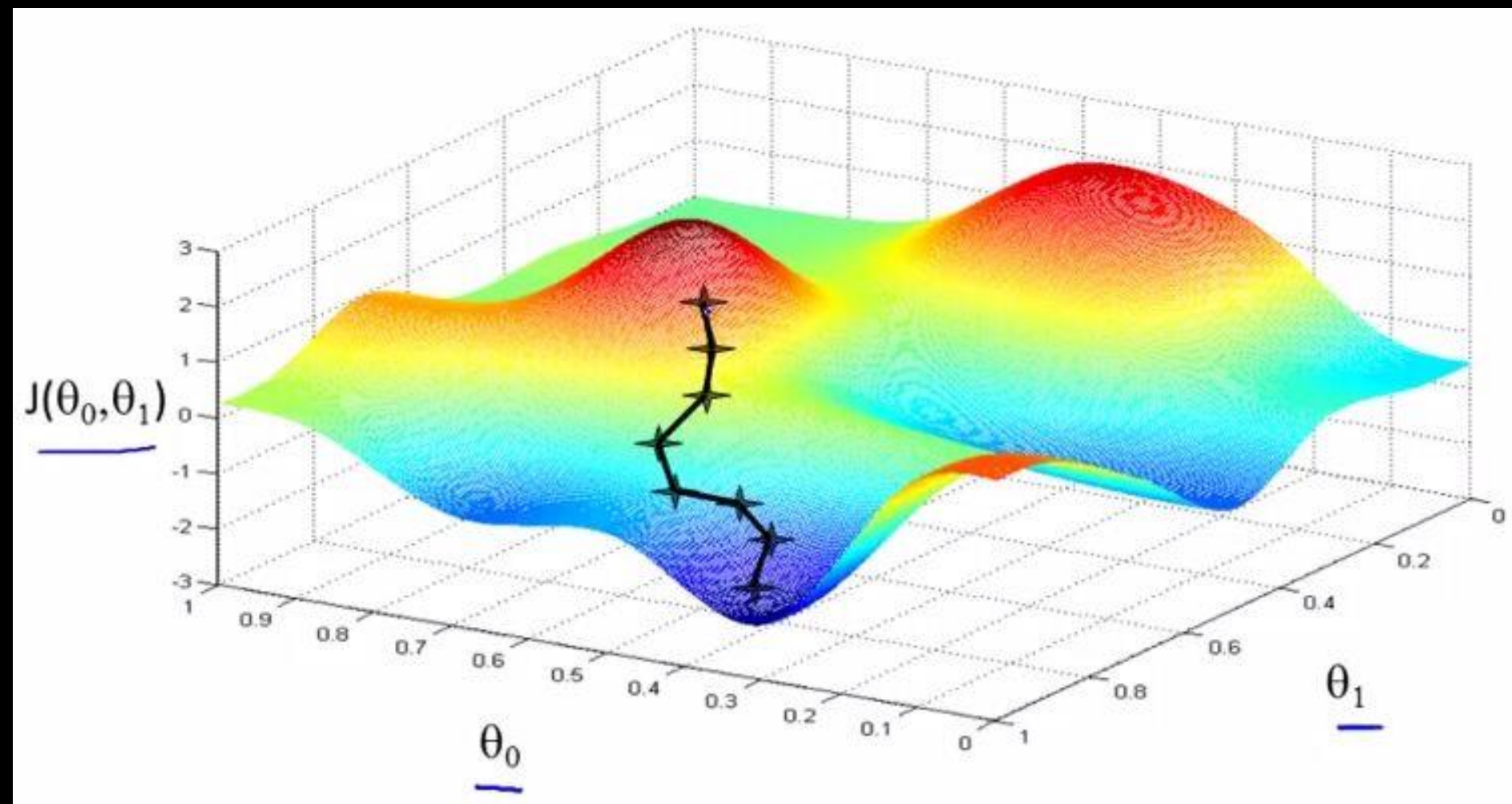
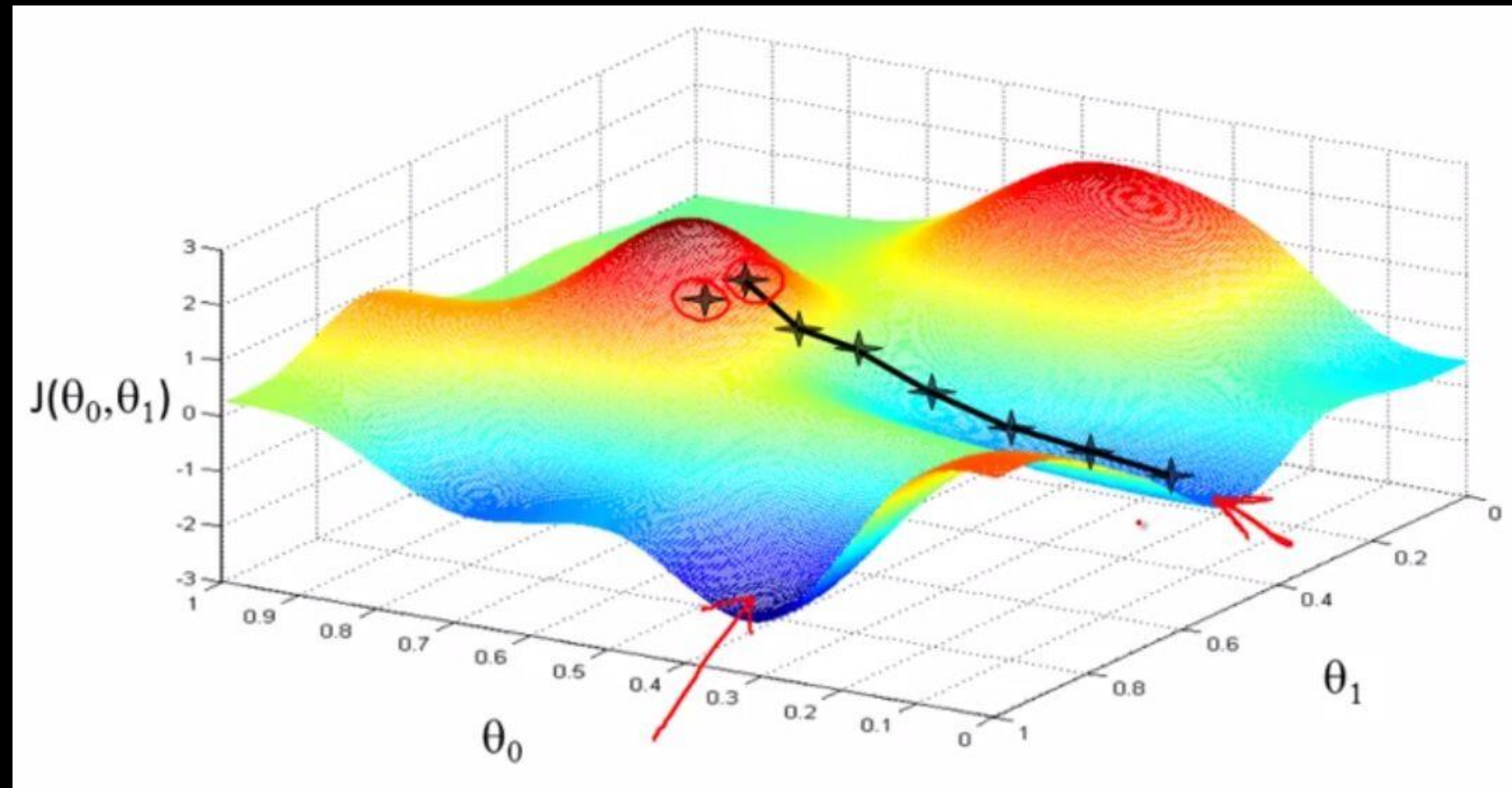
如何找到神經網路的最優化結果

那就沿著能讓誤差下降最快的陡坡那裏走就對了

- 需要是凸函數
- 必需連續可微分



隨機梯度下降/Stochastic Gradient Descent (SGD)



Adam 自我調整動量估計 (Adaptive Moment Estimation)

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$

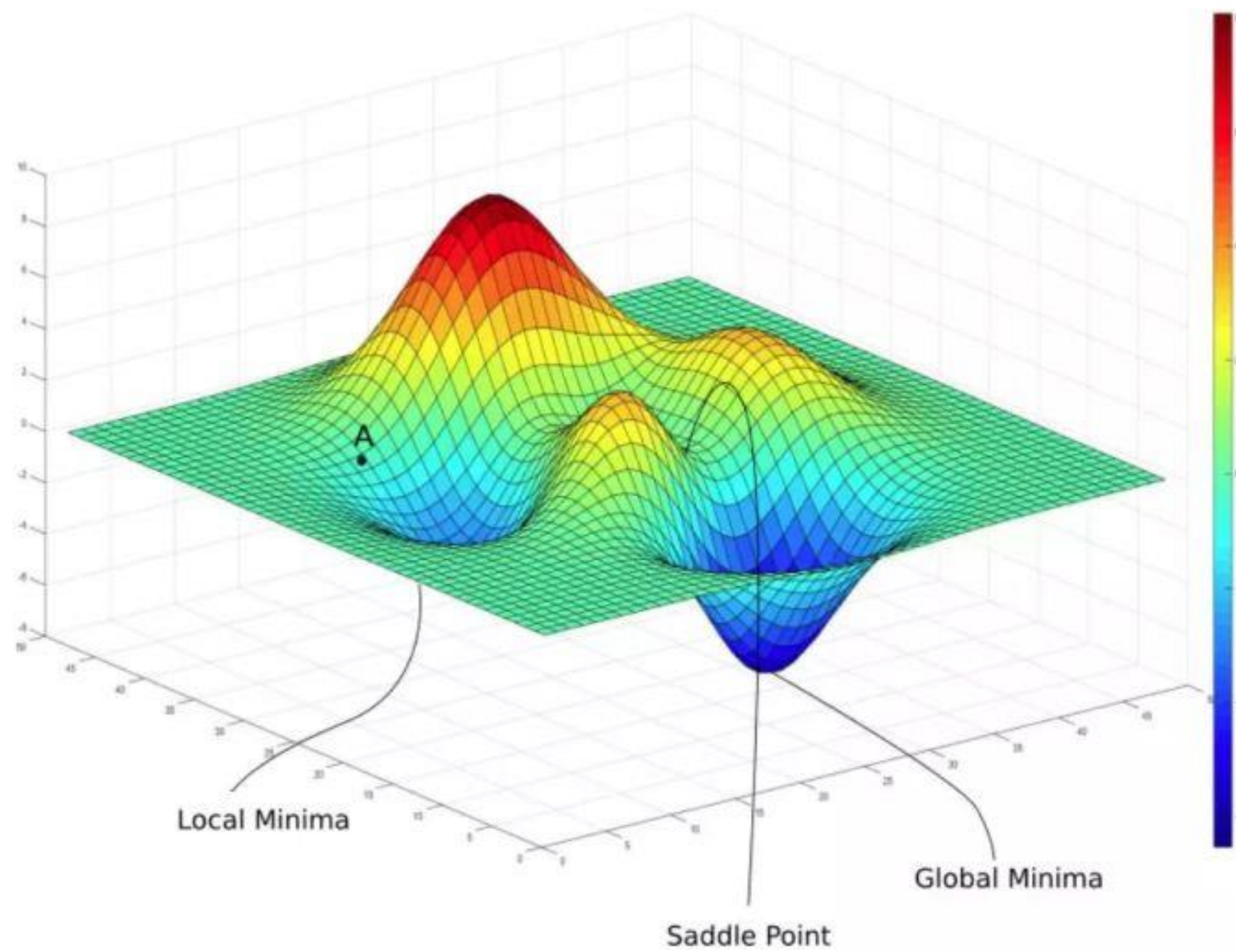
$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

動量

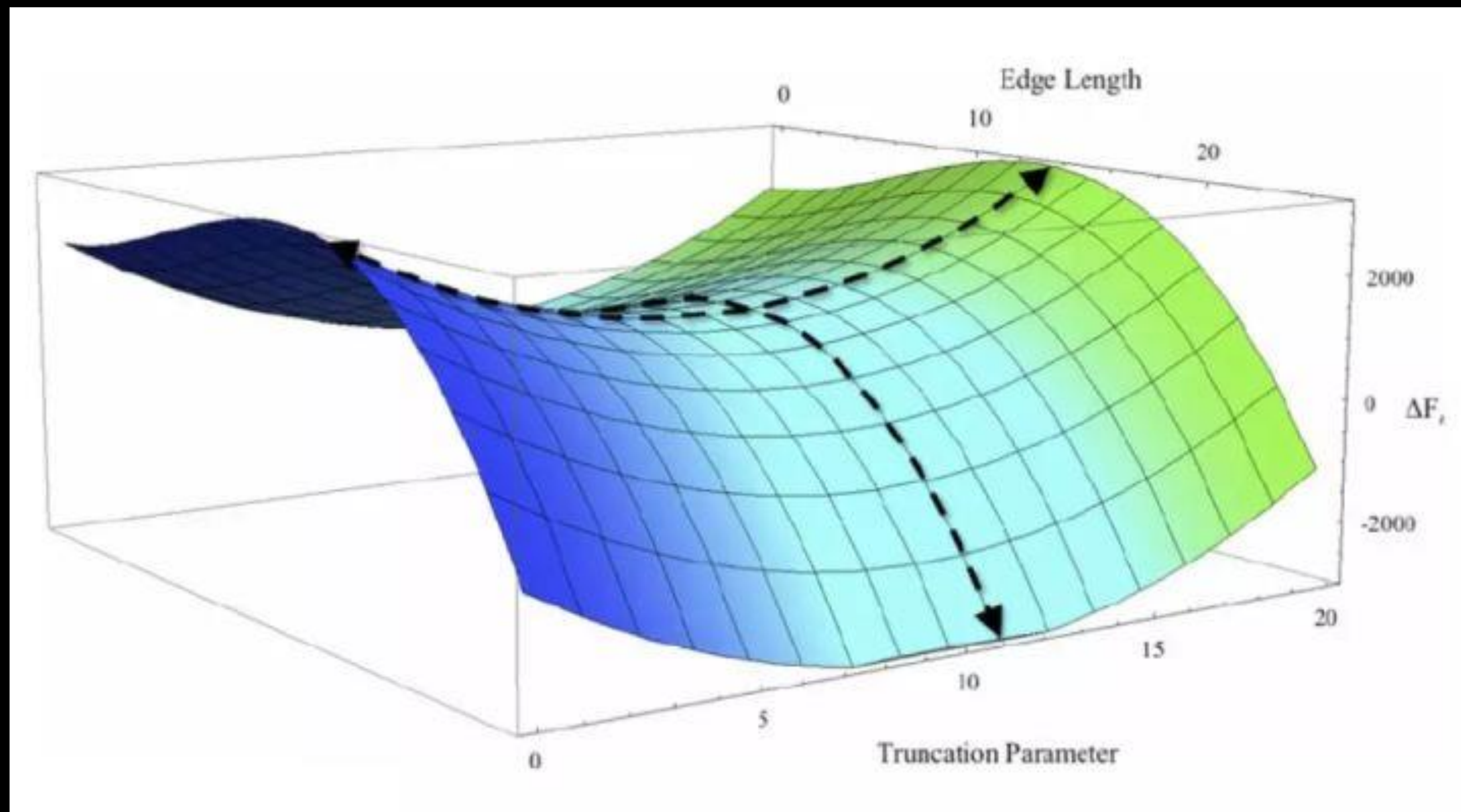
$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

β_1 取0.9, β_2 取0.999, ϵ 取 10^{-8} 。

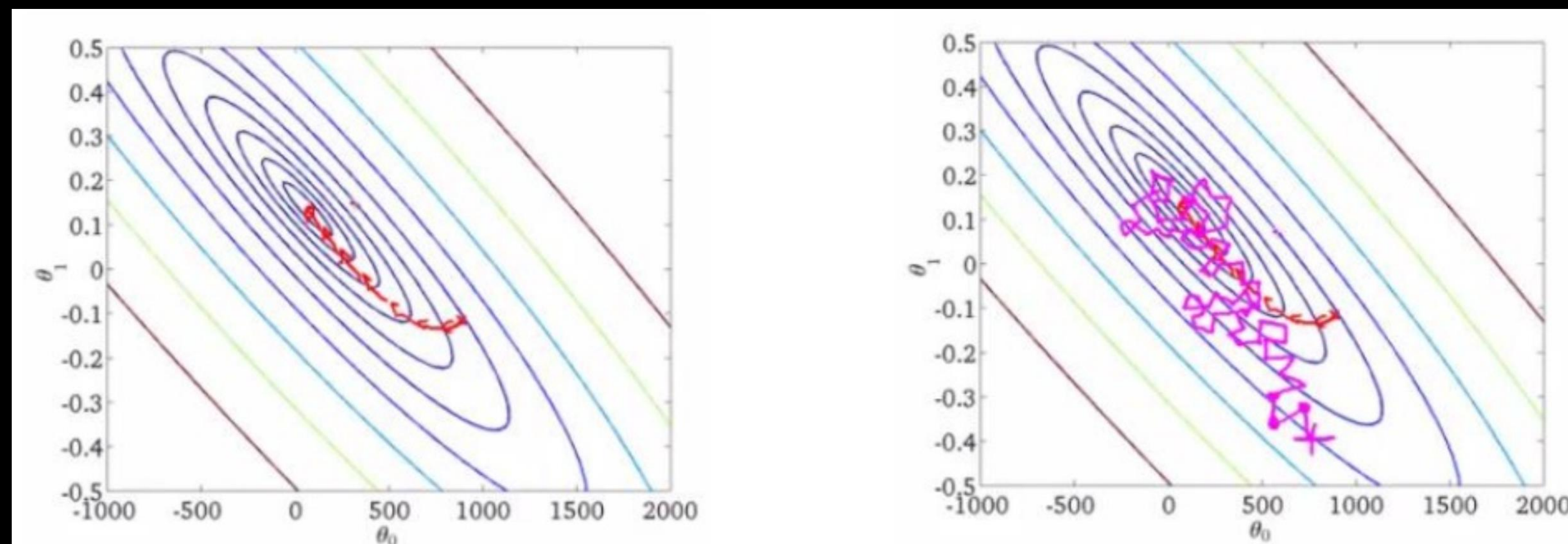
區域最佳解



鞍點



Minibatch



批次梯度下降
緩慢但耗費記憶體

隨機梯度下降
會有波動但快速

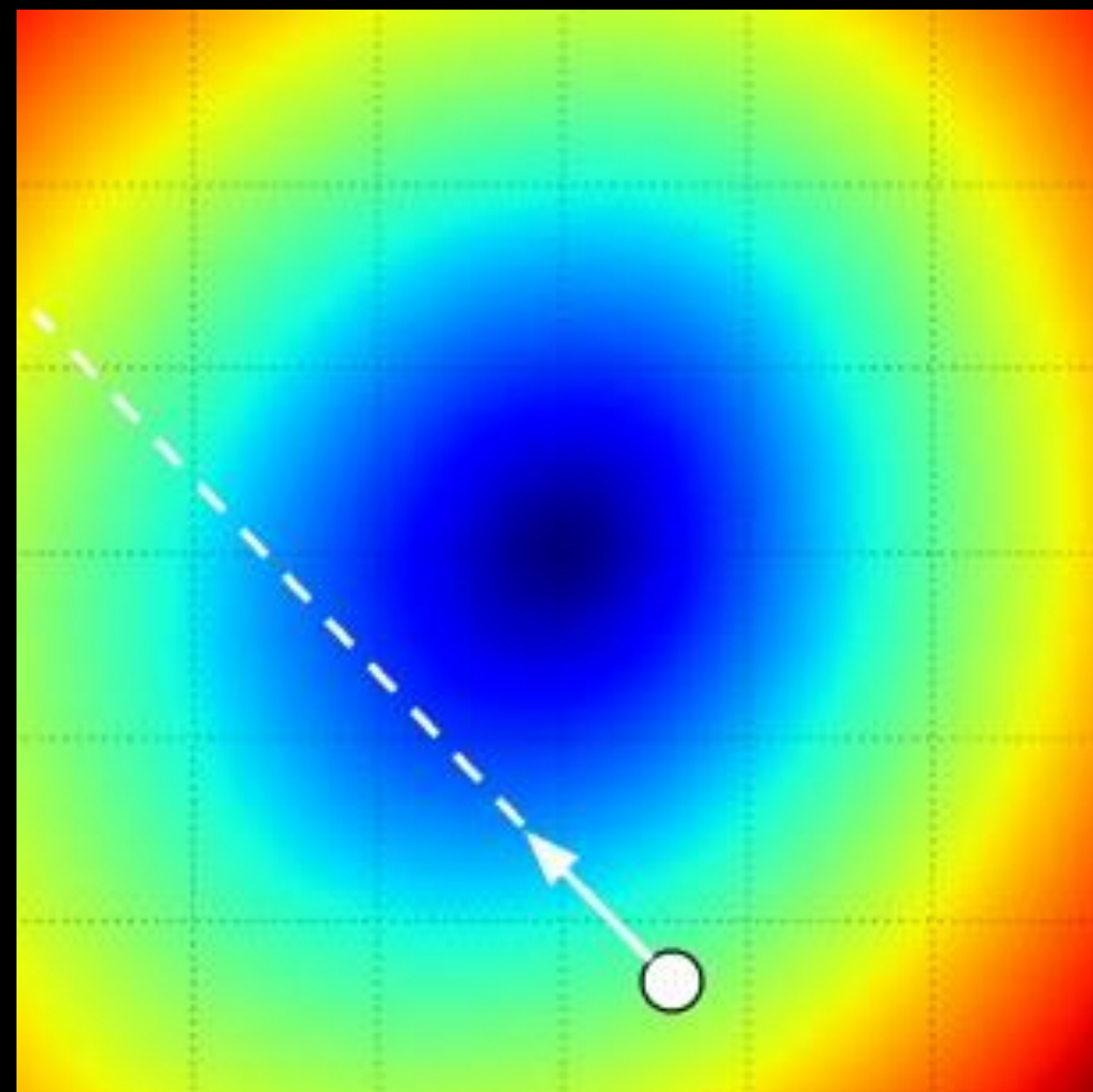
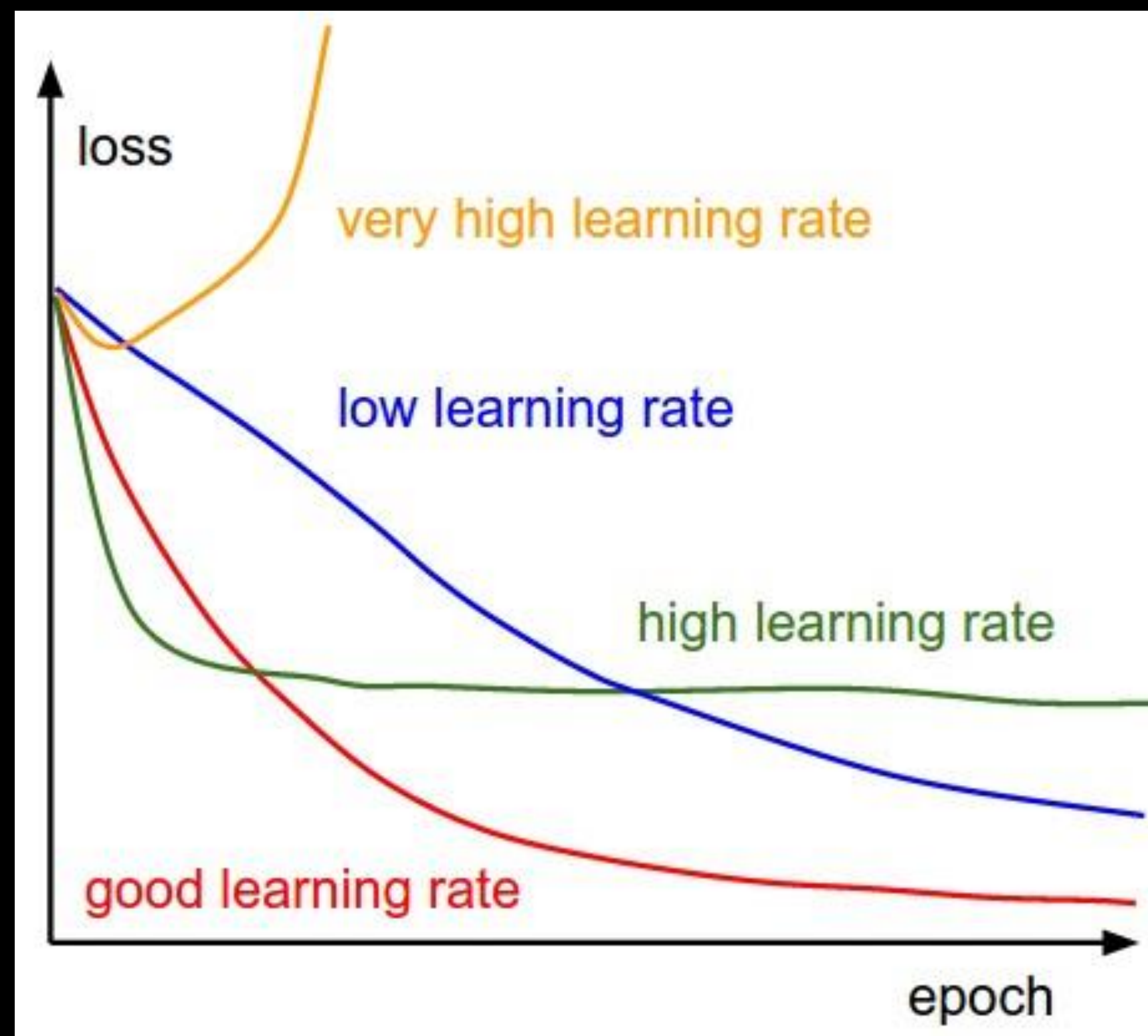


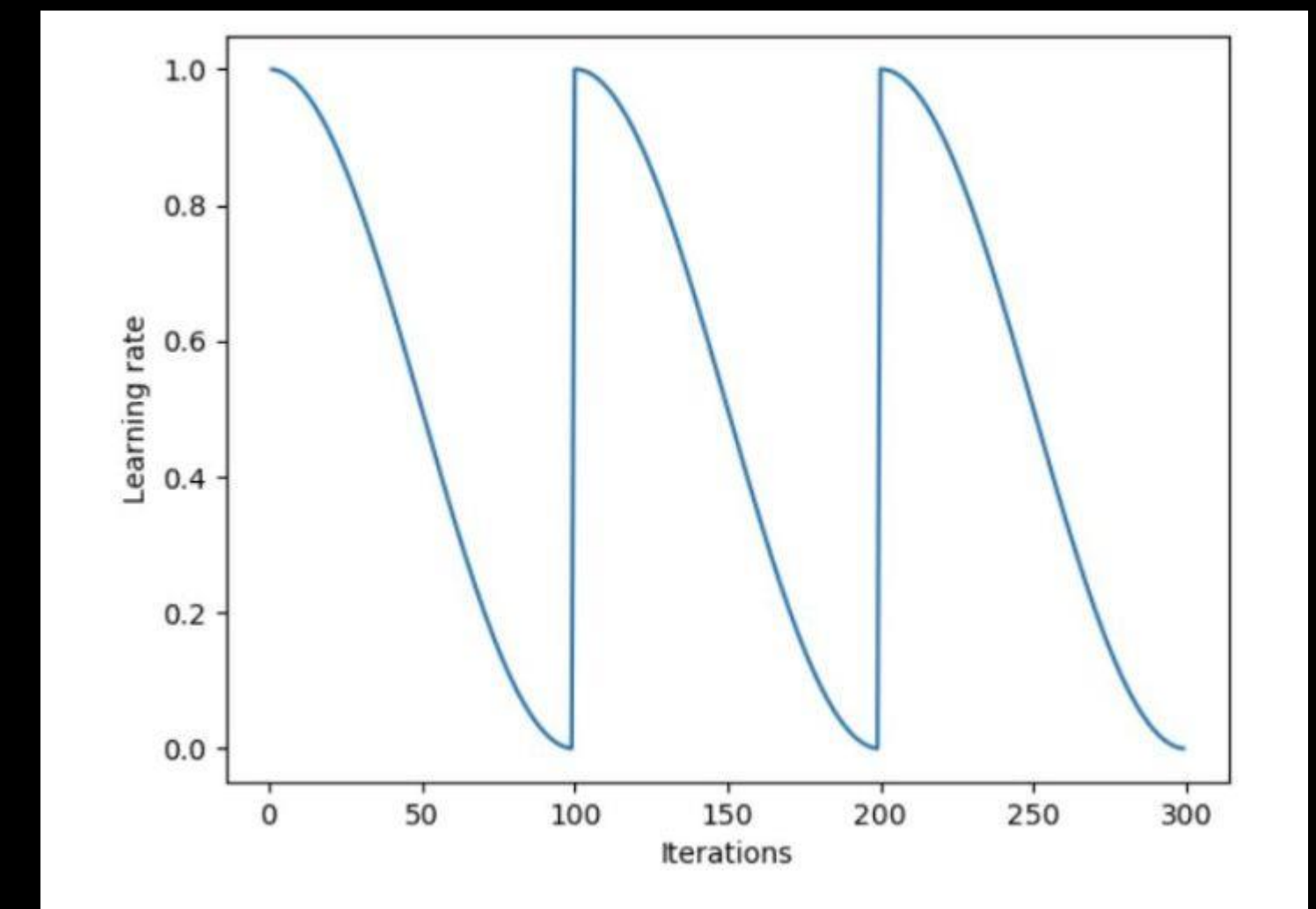
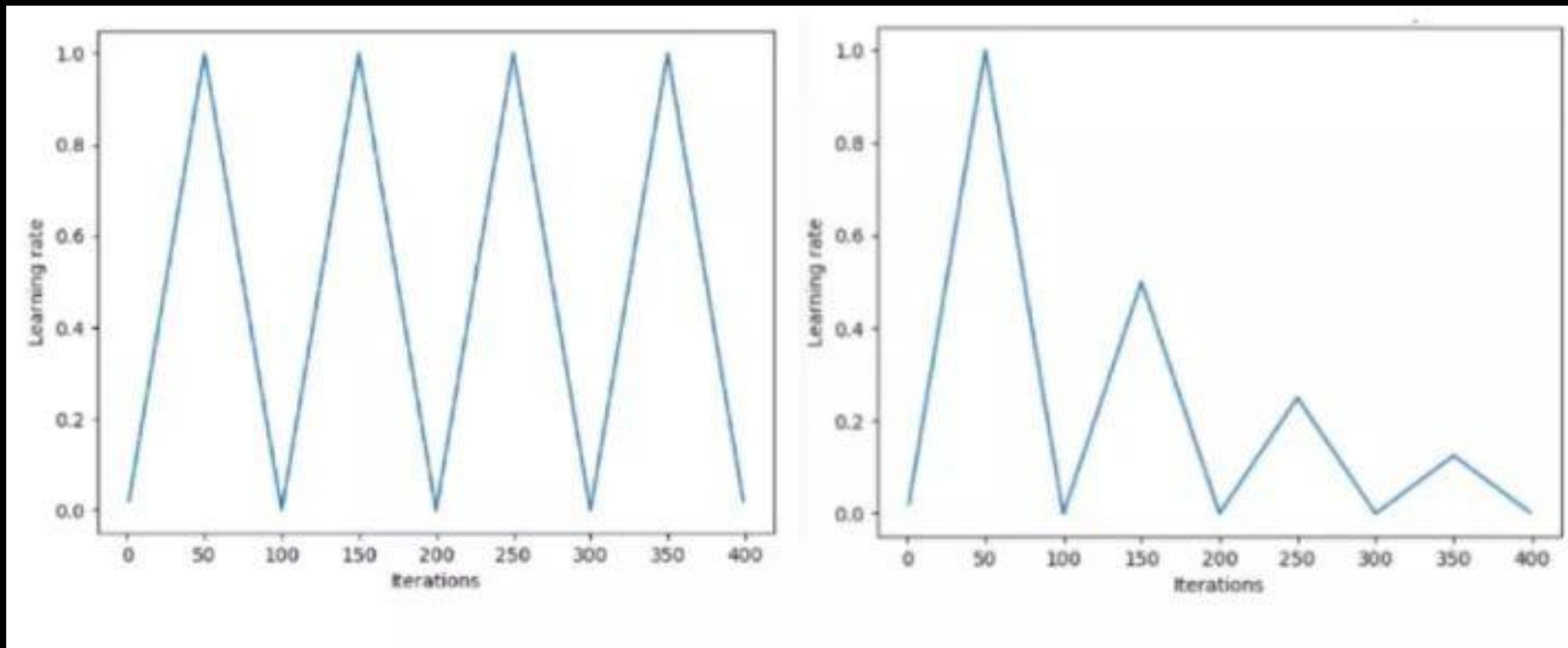
Minibatch更新梯度下降

Epoch

學習速率

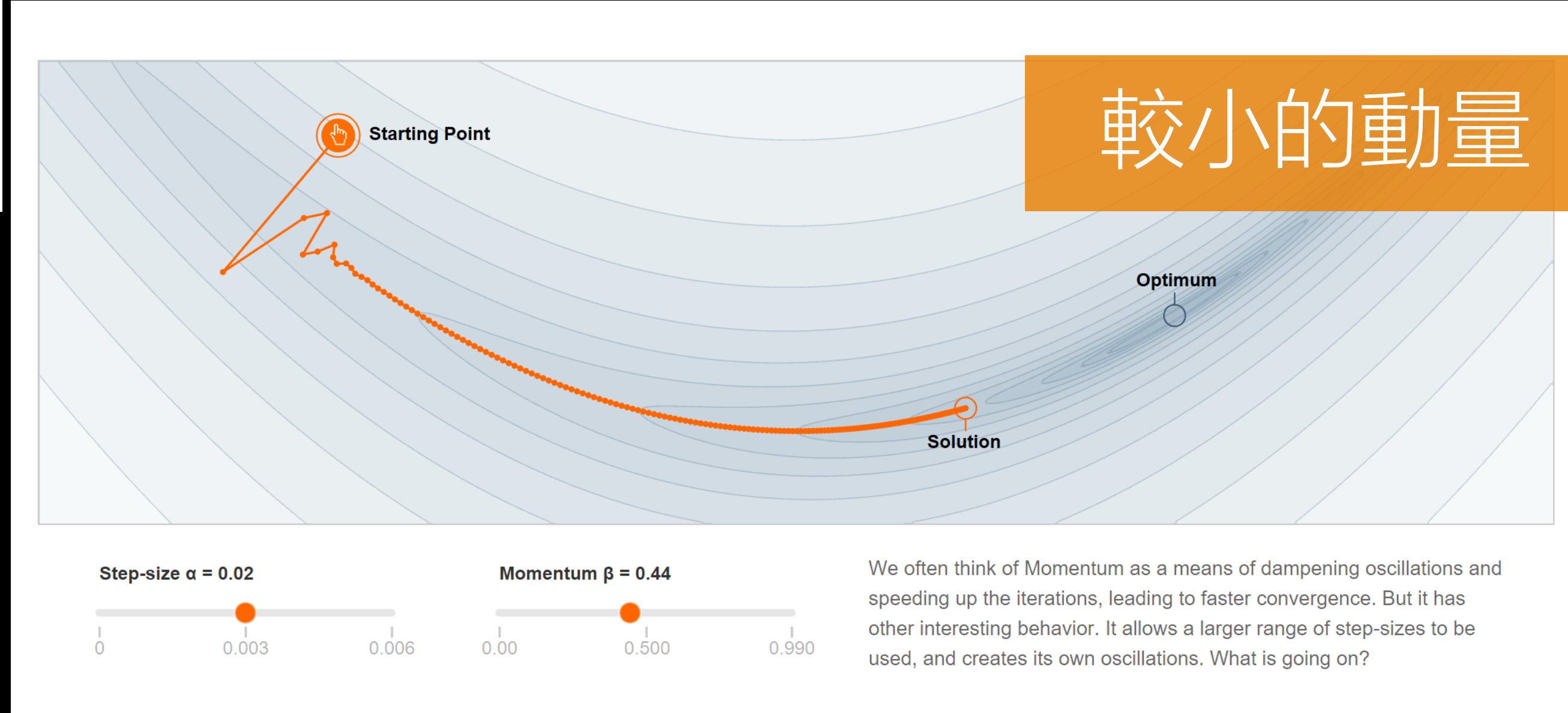
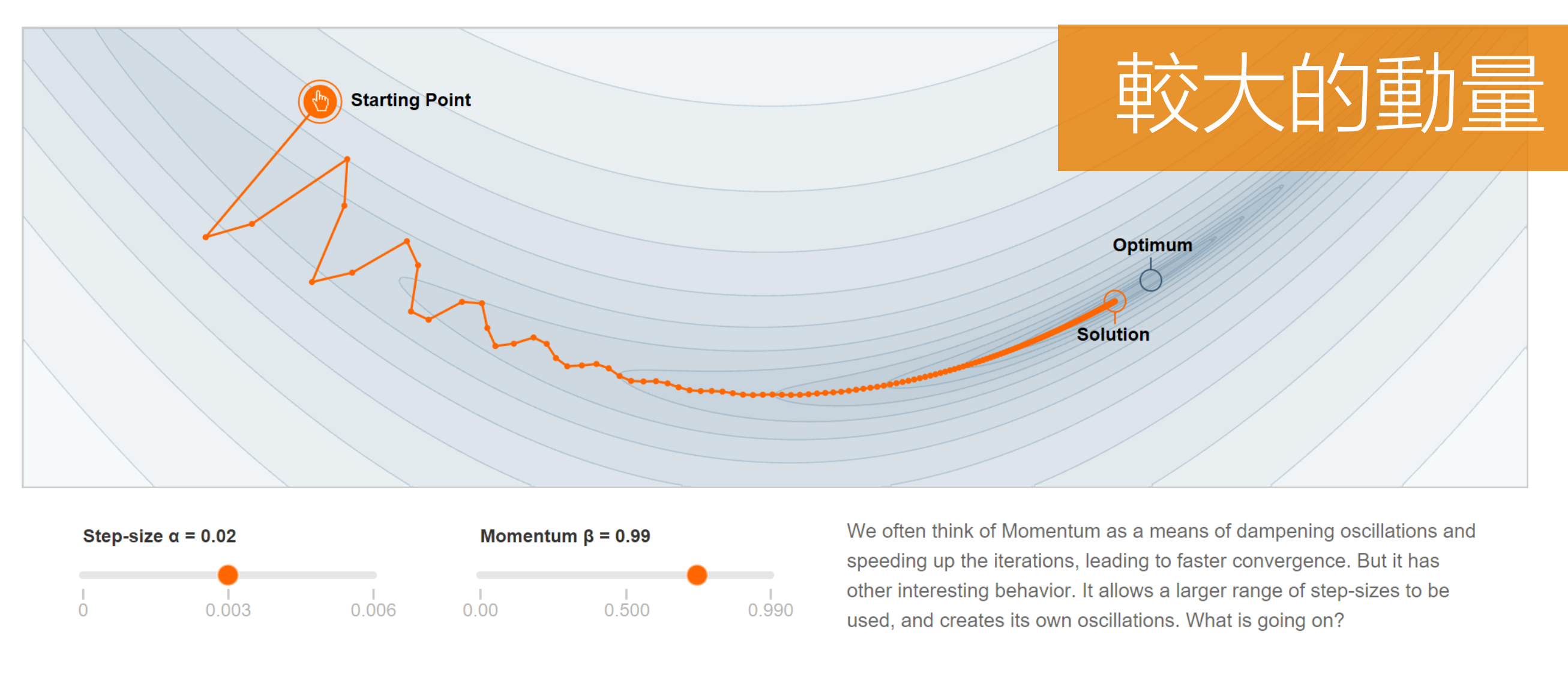
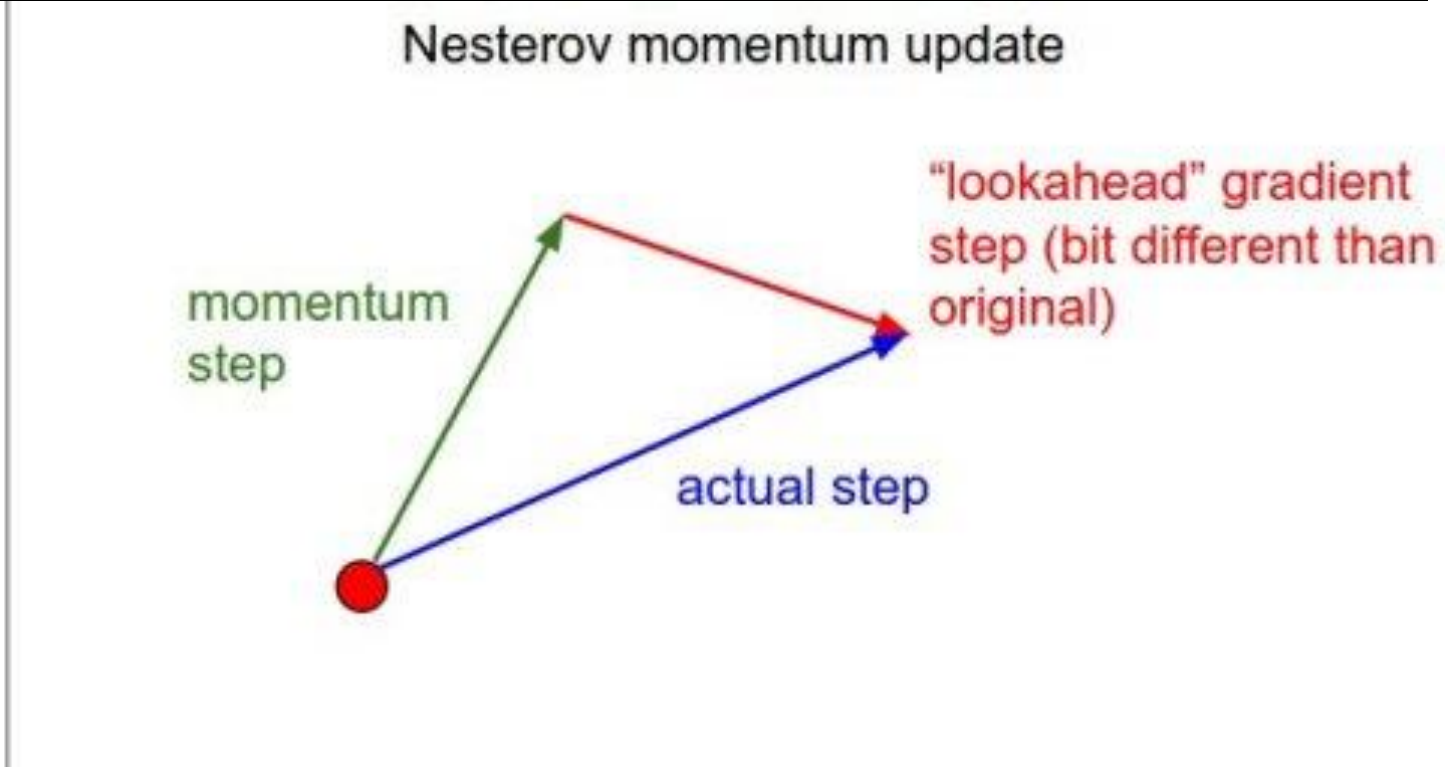
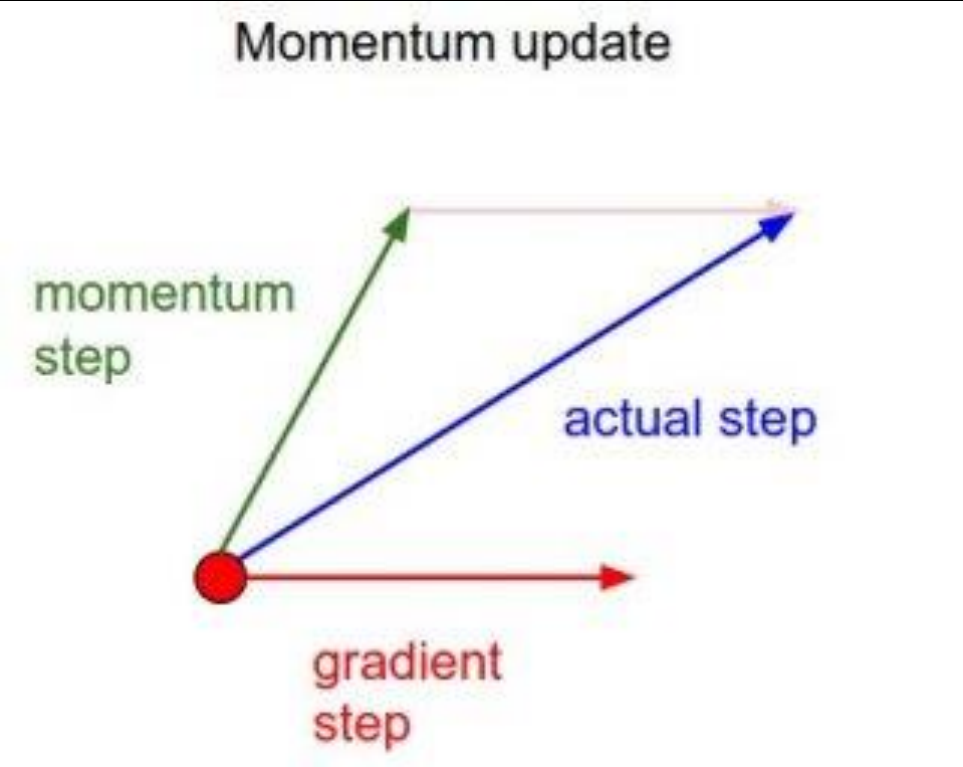
Learning Rate



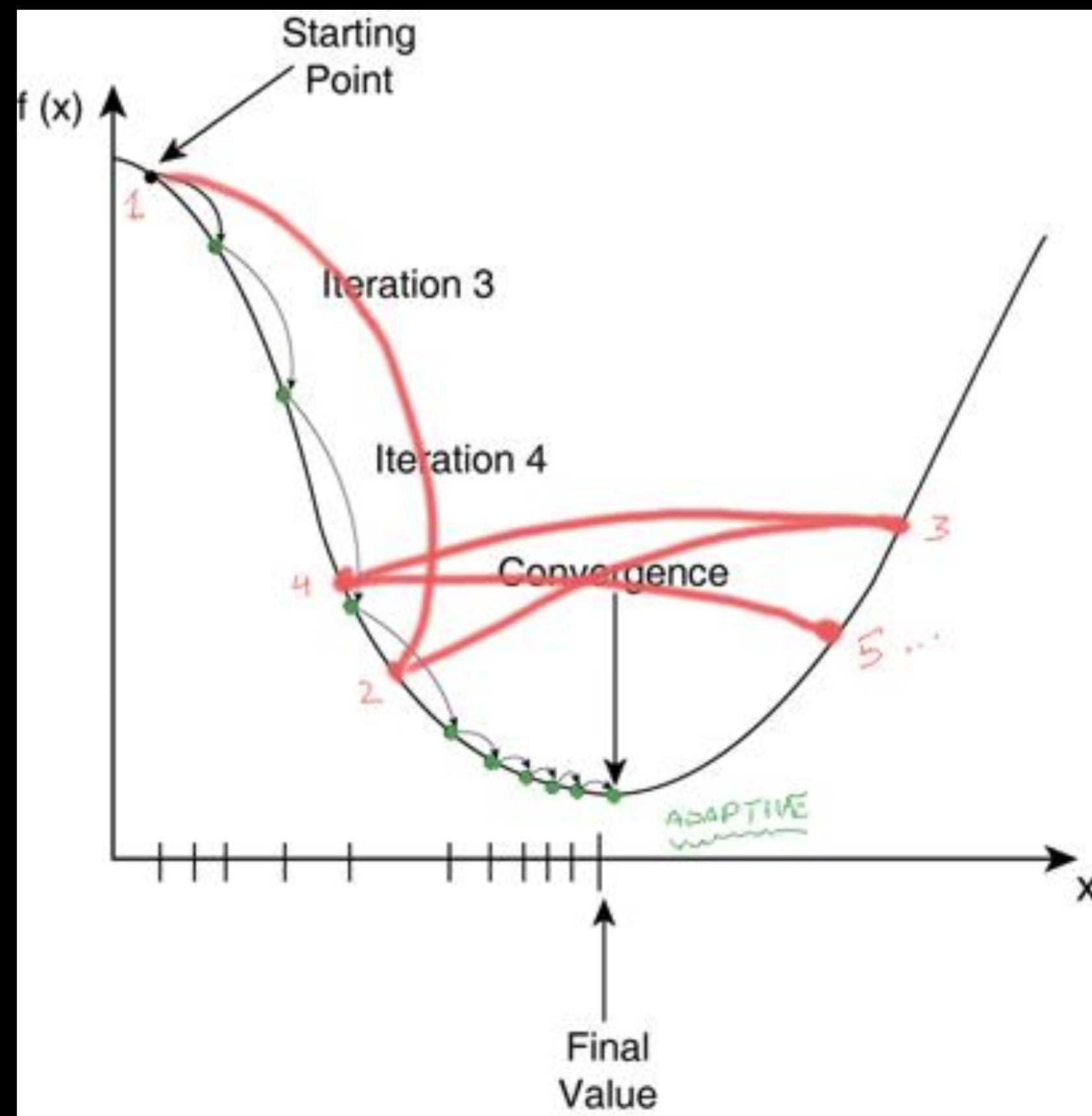


Leslie N. Smith 提出的 Triangular 和 Triangular2 迴圈學習率方法。左側的最大學習率和最小學習率保持不變。右側的區別在於每個週期之後學習率減半

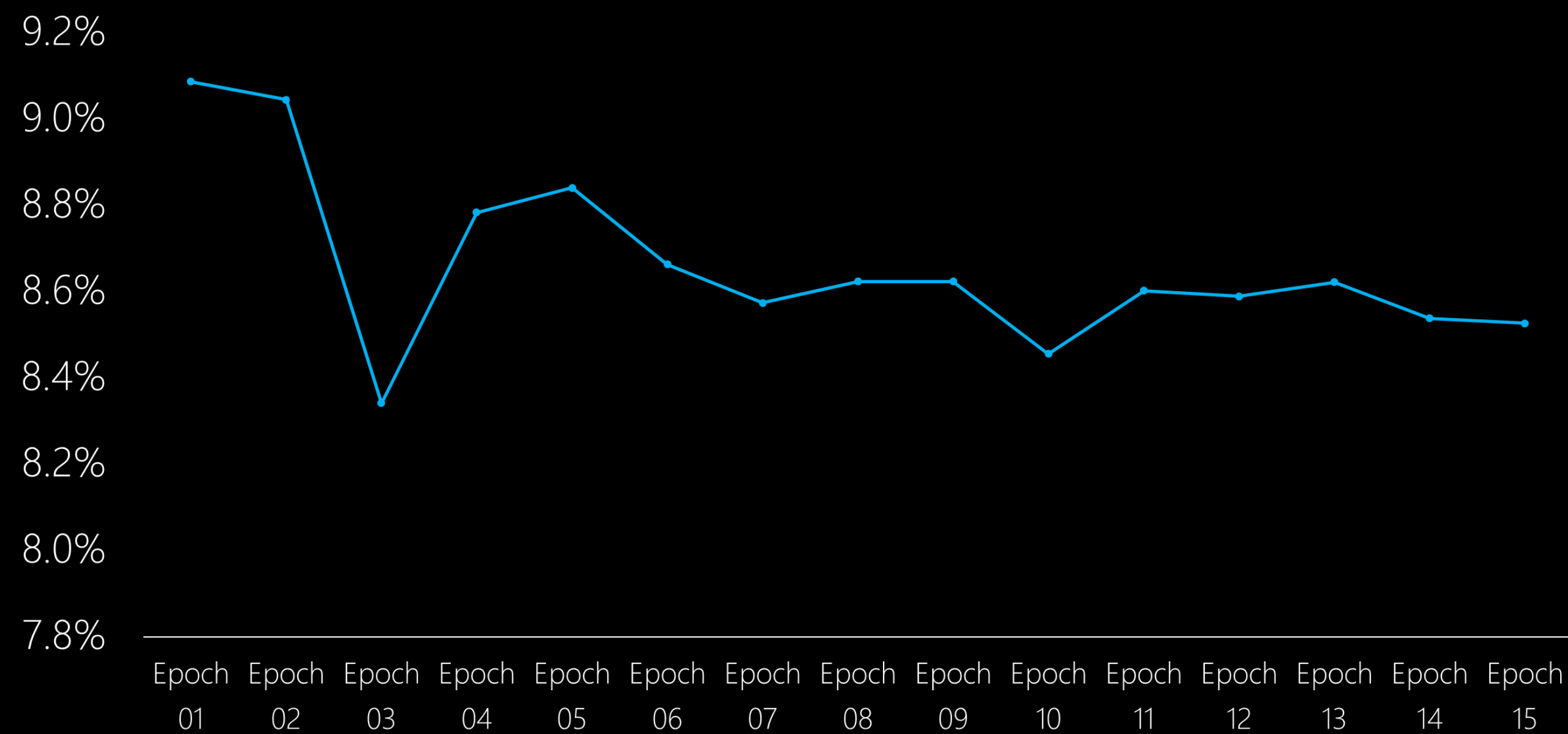
動量 Momentum



這些參數構成了訓練過程

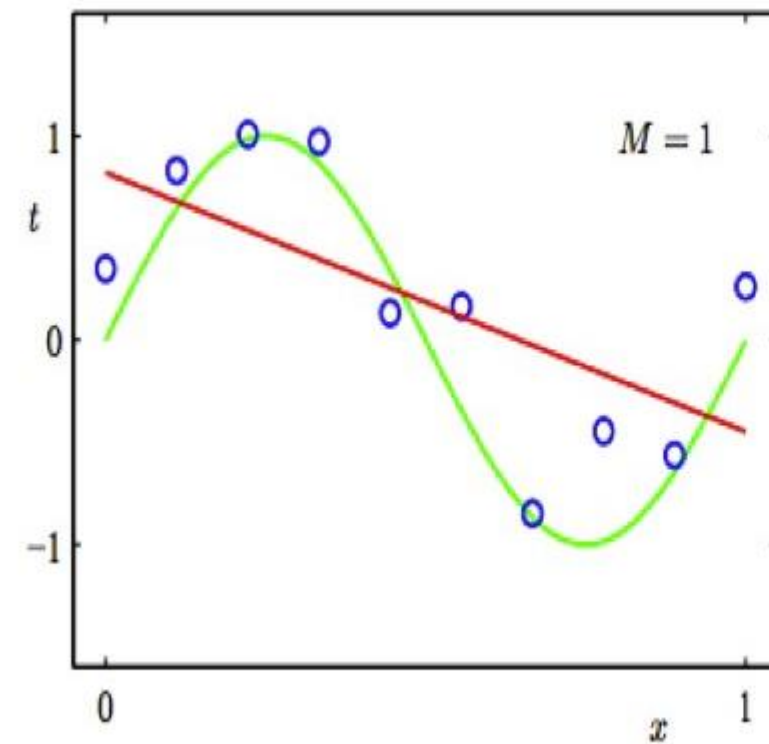


各Epoch的誤差損失收斂

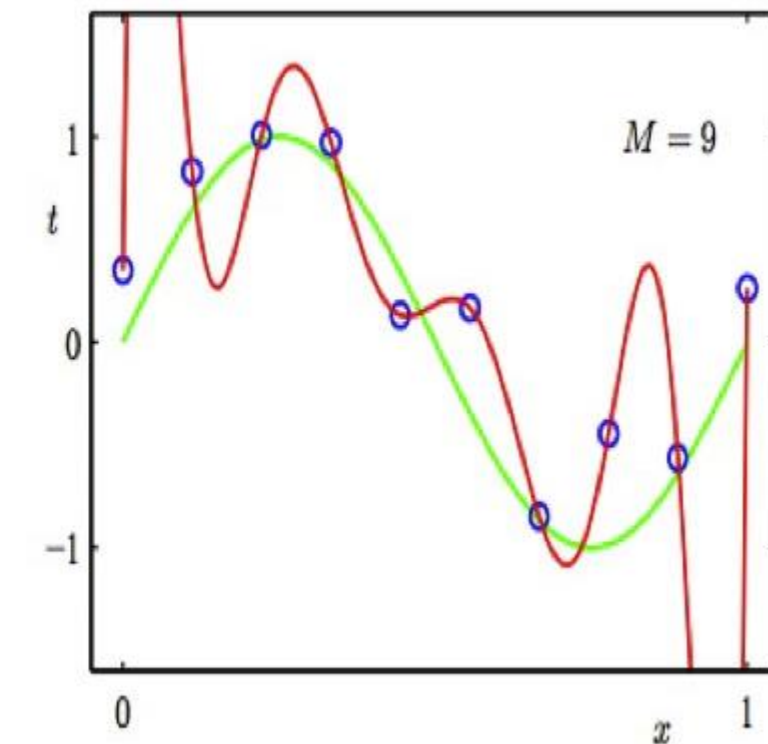
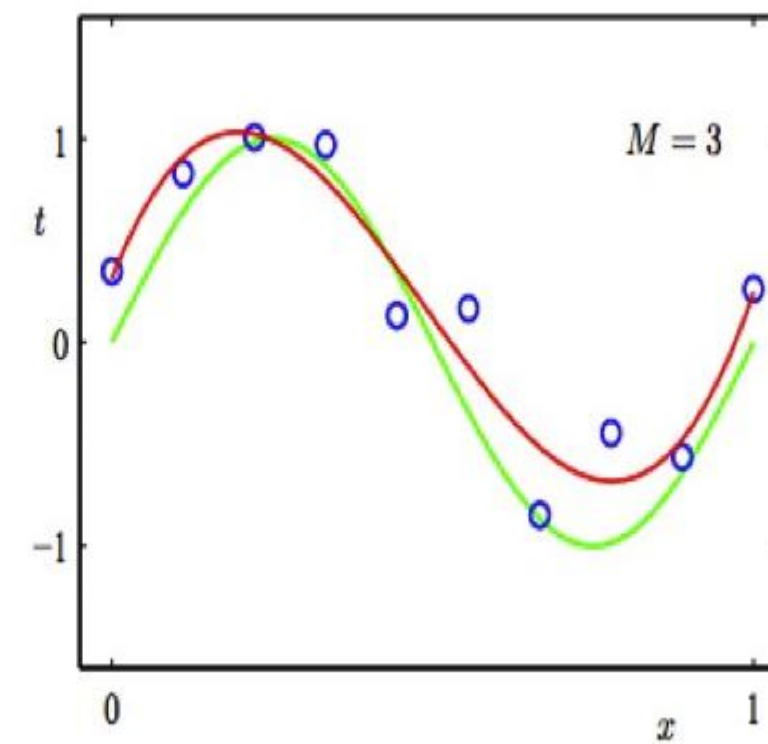


欠擬合(underfitting)與過擬合(Overfitting)

Regression:

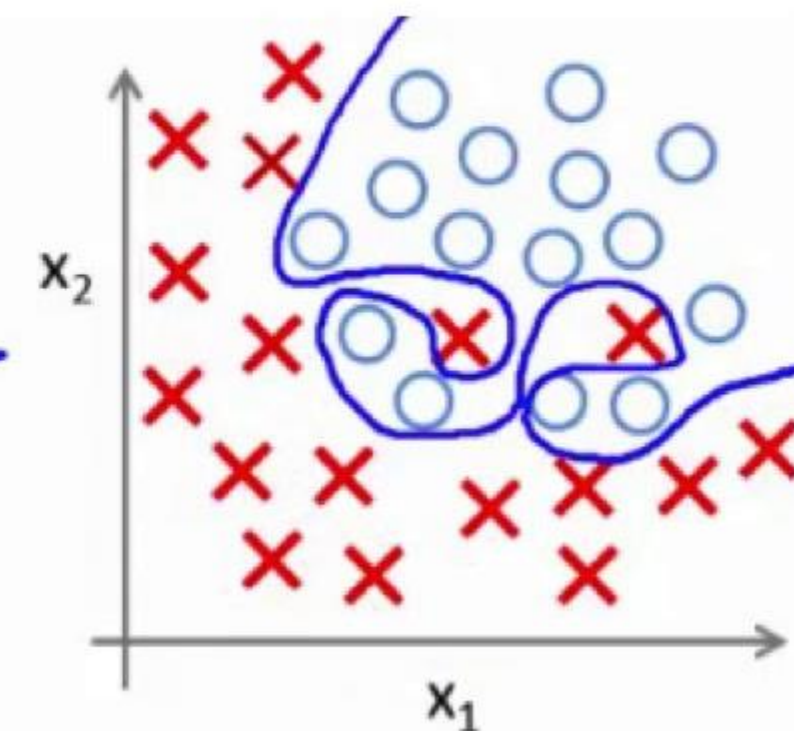
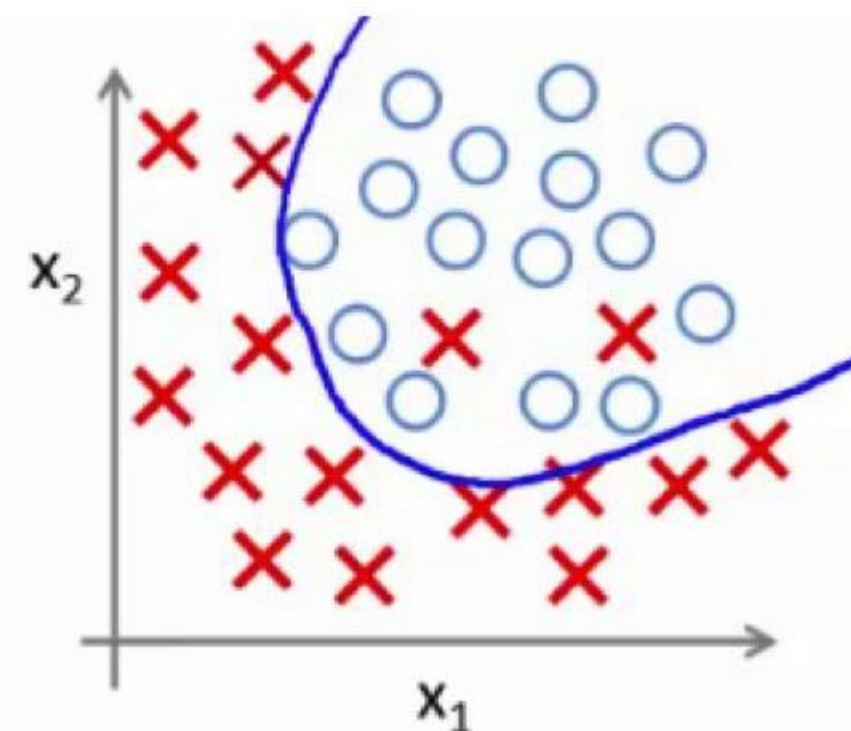
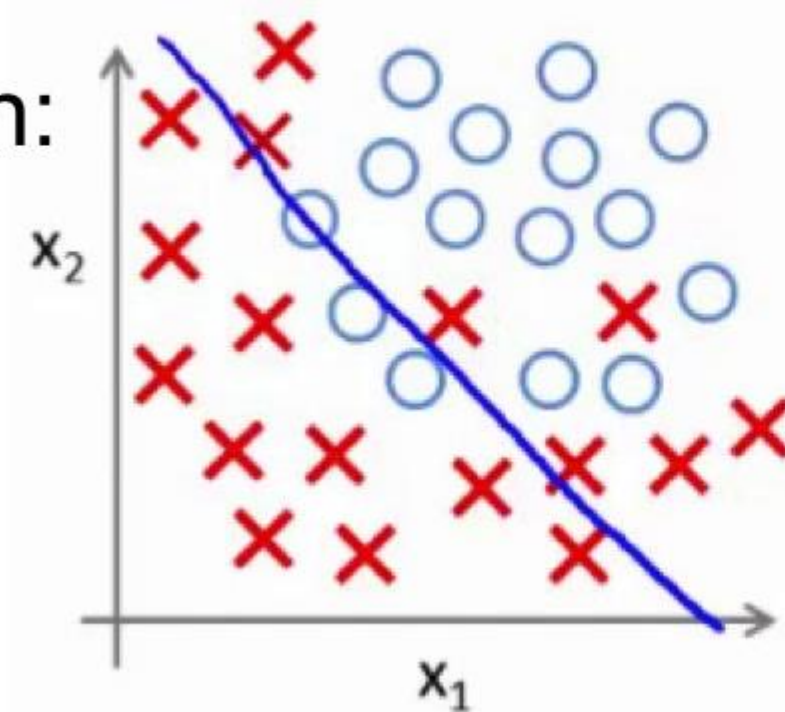


predictor too inflexible:
cannot capture pattern



predictor too flexible:
fits noise in the data

Classification:



L1 regulariZation L1 正則會讓模型變稀疏

$$C = C_0 + \frac{\lambda}{n} \sum_w |w|.$$

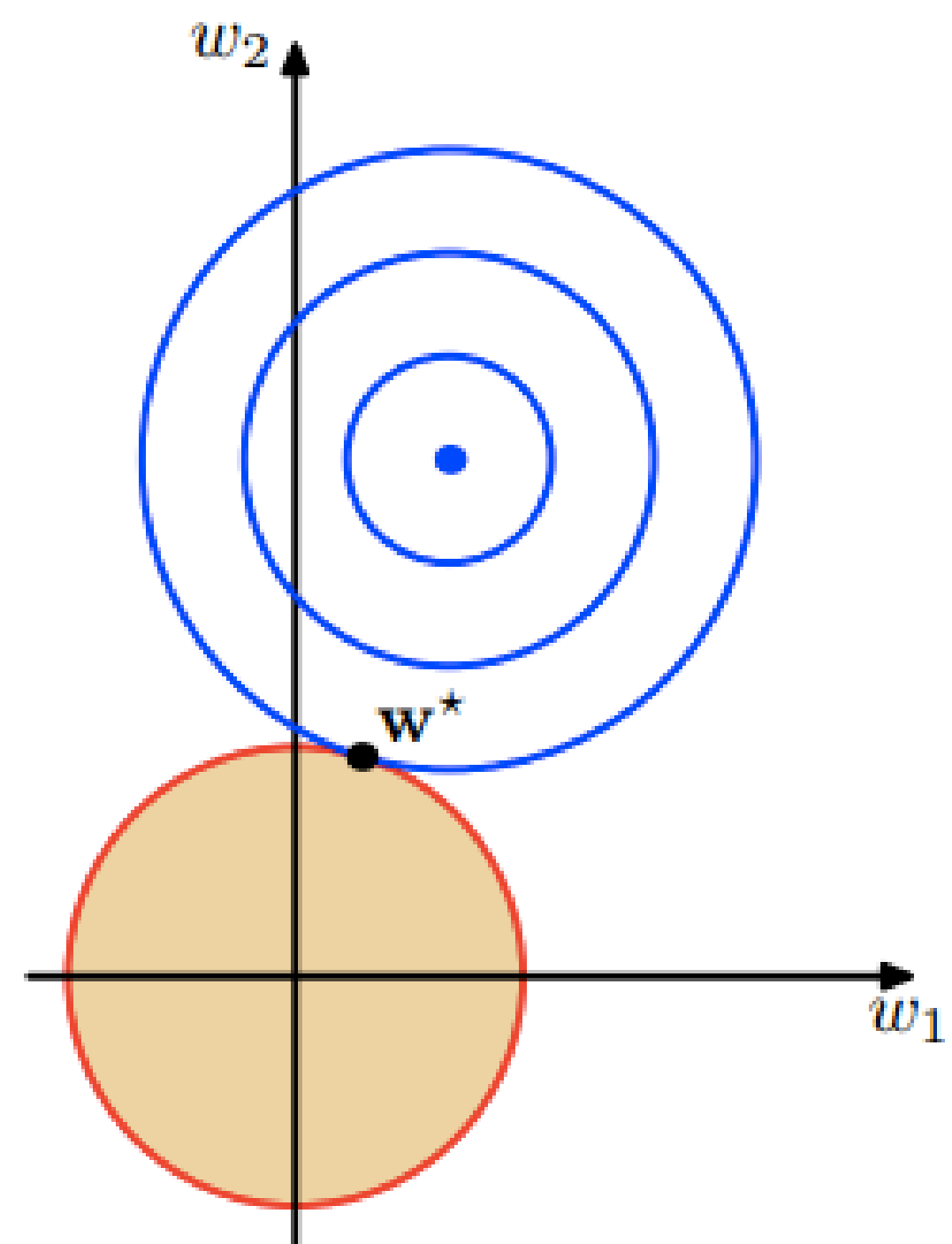
因為有了L1正則的懲罰項，
因此傾向往部分權重靠攏，
其餘變零

L2 regularization L2正則又稱之為權重遞減

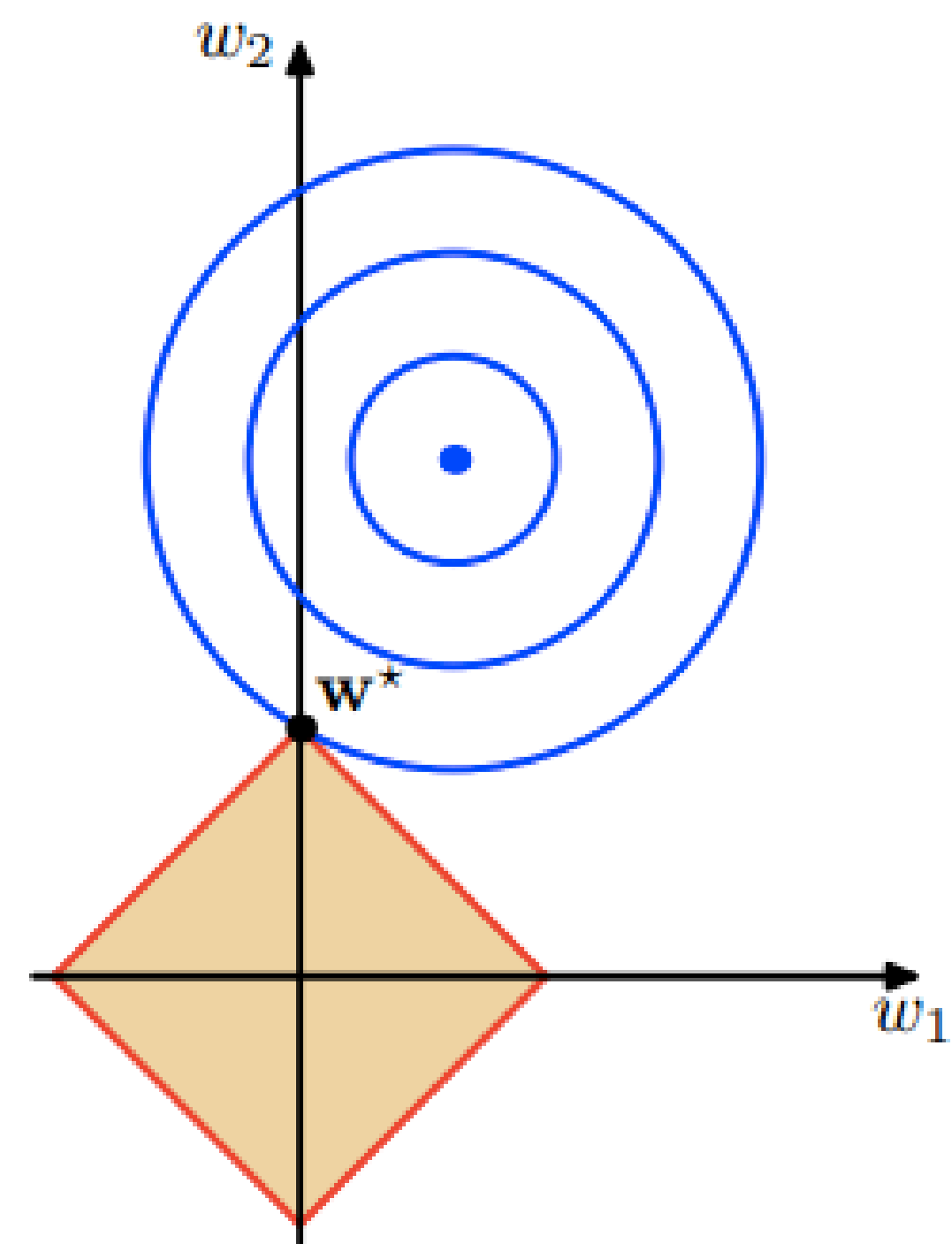
$$C = C_0 + \frac{\lambda}{2n} \sum_w w^2,$$

因為有了L2正則的懲罰項，
因此傾向往全體權重最小的
方向邁進

L2正則



L1正則



那些因素會影響模型的成果

數據

模型結構

超參數


避免過擬合

最佳化方法

正歸化

活化函數

權重初始化



Q & A