

類神經網路訓練不起來怎麼辦 (一)：局部最小值 (local minima) 與鞍點 (saddle point)

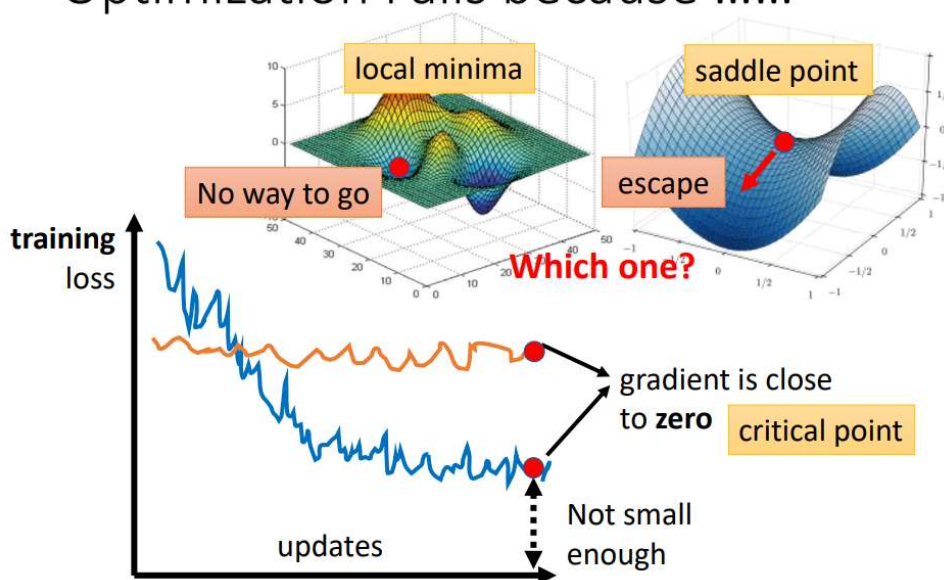
Create at 2022/06/08

- 類神經網路訓練不起來怎麼辦 (一)：局部最小值 (local minima) 與鞍點 (saddle point).
 - Optimization 時怎麼把 gradient descent 做得更好
 - 為甚麼 optimization 會失敗呢？
 - Saddle Point 跟 Local Minima 誰比較常見？
- 上課資源：
 1. 類神經網路訓練不起來怎麼辦 (一)：局部最小值 (local minima) 與鞍點 (saddle point).
(<https://www.youtube.com/watch?v=QW6uINn7uGk>).

Optimization 時怎麼把 gradient descent 做得更好

為甚麼 optimization 會失敗呢？

Optimization Fails because



- 隨著參數不斷 update，training loss 不會再下降
- 但是對於這個 loss 仍然不滿意
- 有時候會發現 model train 不起來，不管怎麼 update 參數 loss 都掉不下去
 - 猜想：走到一個點，這個點對 loss 的微分為 0，當對 loss 的微分為 0 時，gradient descent 就沒辦法再 update 參數，此時 training 就停下來了，參數不再 update，loss 就不會再下降了
- 不是只有 local minima 的 gradient 為 0
- saddle point 的 gradient 是 0，但不是 local minima 也不是 local maxima
- gradient 為 0 的點，統稱為 critical point
- 可以說 loss 沒辦法再下降，可能是卡在 critical point

Warning of Math

可以跳過沒關係

Taylor Series Approximation

$L(\theta)$ around $\theta = \theta'$ can be approximated below

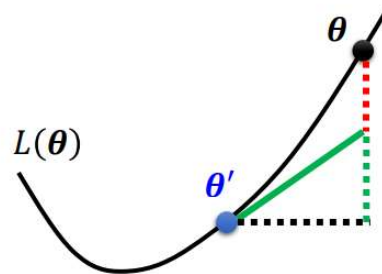
$$L(\theta) \approx L(\theta') + (\theta - \theta')^T \mathbf{g} + \frac{1}{2} (\theta - \theta')^T \mathbf{H} (\theta - \theta')$$

Gradient \mathbf{g} is a vector

$$\mathbf{g} = \nabla L(\theta') \quad g_i = \frac{\partial L(\theta')}{\partial \theta_i}$$

Hessian \mathbf{H} is a matrix

$$H_{ij} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} L(\theta')$$



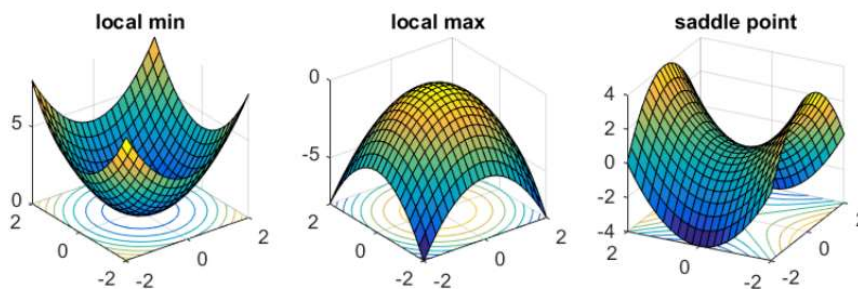
Hessian

$L(\theta)$ around $\theta = \theta'$ can be approximated below

$$L(\theta) \approx L(\theta') + (\theta - \theta')^T \mathbf{g} + \frac{1}{2} (\theta - \theta')^T \mathbf{H} (\theta - \theta')$$

At critical point

telling the properties of critical points



Hessian

At critical point: $v^T H v$

$$L(\theta) \approx L(\theta') + \frac{1}{2}(\theta - \theta')^T H(\theta - \theta')$$

For all v

$$v^T H v > 0 \implies \text{Around } \theta': L(\theta) > L(\theta') \implies \text{Local minima}$$

= H is positive definite = All eigen values are positive. \uparrow

For all v

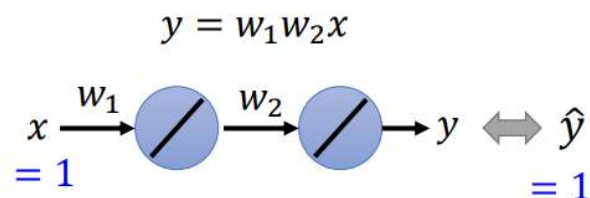
$$v^T H v < 0 \implies \text{Around } \theta': L(\theta) < L(\theta') \implies \text{Local maxima}$$

= H is negative definite = All eigen values are negative. \uparrow

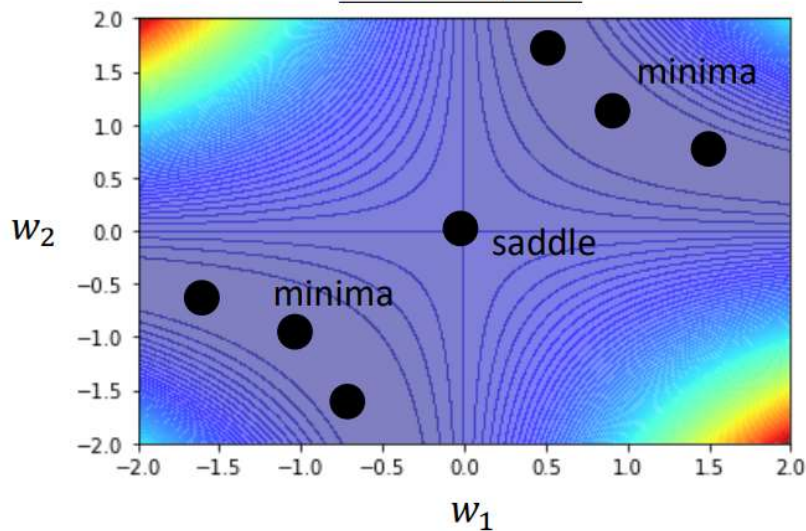
Sometimes $v^T H v > 0$, sometimes $v^T H v < 0 \implies \text{Saddle point}$

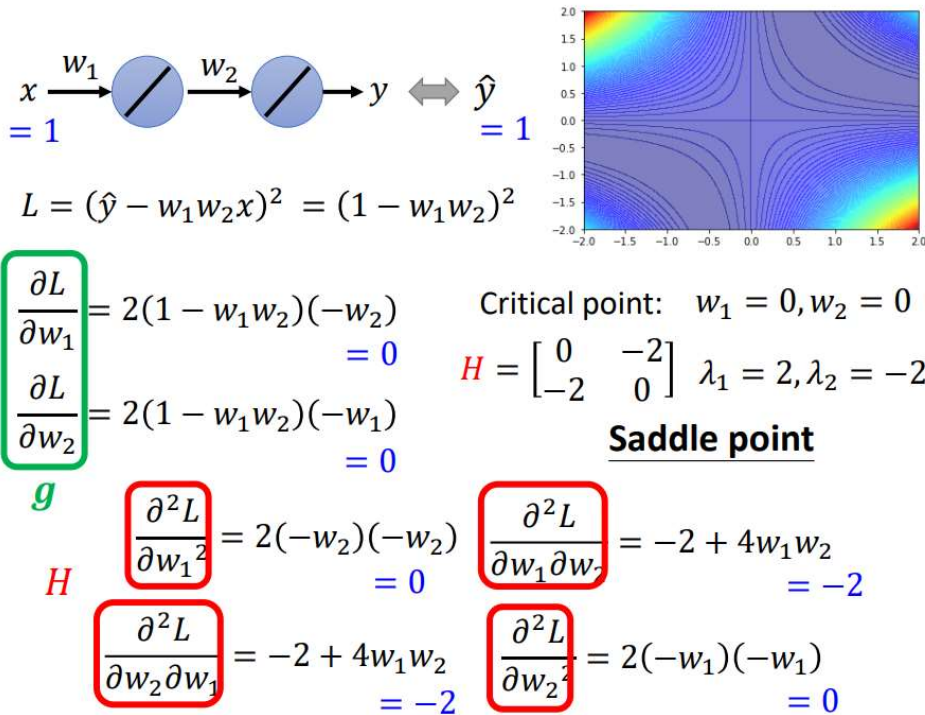
Some eigen values are positive, and some are negative. \uparrow

Example



Error Surface





Don't afraid of saddle point?

$$v^T H v$$

At critical point: $L(\theta) \approx L(\theta') + \frac{1}{2}(\theta - \theta')^T H(\theta - \theta')$

Sometimes $v^T H v > 0$, sometimes $v^T H v < 0 \Rightarrow$ Saddle point

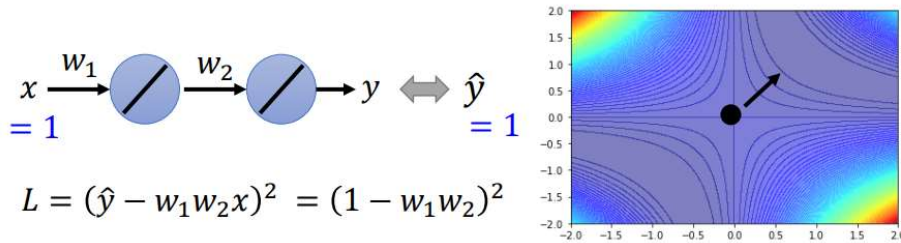
H may tell us parameter update direction!

u is an eigen vector of H
 λ is the eigen value of u
 $\lambda < 0$

$$u^T H u = u^T (\lambda u) = \lambda \|u\|^2 < 0$$

$$L(\theta) \approx L(\theta') + \frac{1}{2}(\theta - \theta')^T H(\theta - \theta') \Rightarrow L(\theta) < L(\theta')$$

$$\theta - \theta' = u \quad \theta = \theta' + u \quad \text{Decrease } L$$



$$\frac{\partial L}{\partial w_1} = 2(1 - w_1 w_2)(-w_2)$$

Critical point: $w_1 = 0, w_2 = 0$

$$\frac{\partial L}{\partial w_2} = 2(1 - w_1 w_2)(-w_1)$$

$$H = \begin{bmatrix} 0 & -2 \\ -2 & 0 \end{bmatrix} \quad \lambda_1 = 2, \lambda_2 = -2$$

Saddle point

$\lambda_2 = -2$ Has eigenvector $\mathbf{u} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

Update the parameter along the direction of \mathbf{u}

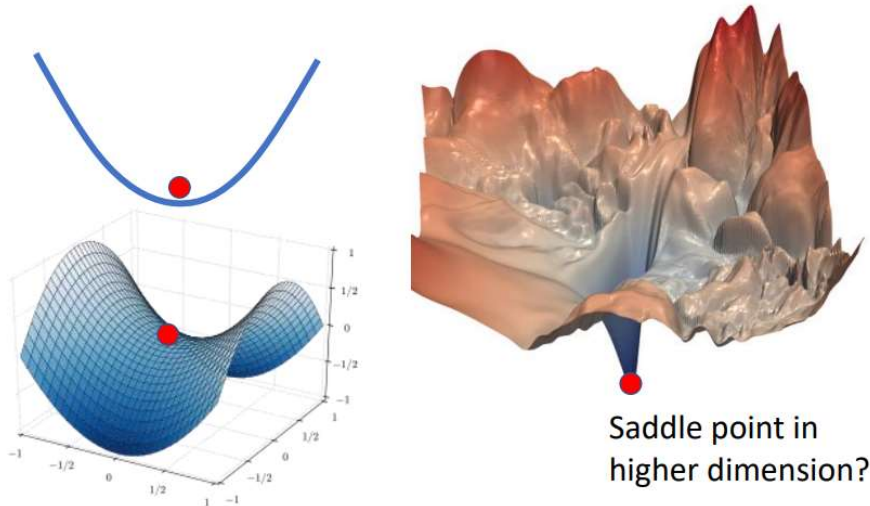
You can escape the saddle point and decrease the loss.

(this method is seldom used in practice)

End of Warning

Saddle Point 跟 Local Minima 誰比較常見？

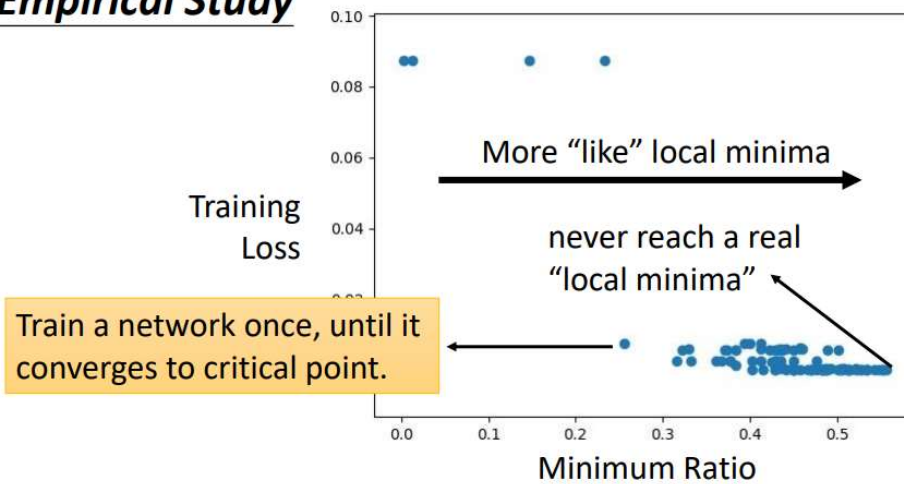
Saddle Point v.s. Local Minima



When you have lots of parameters, perhaps local minima is rare?

- 維度越高，可能可以走的路越多
- 所以在訓練 network 時，參數往往會很多，所以 error surface 其實是在一個非常高的維度中
- 參數有多少代表 error surface 的維度有多少

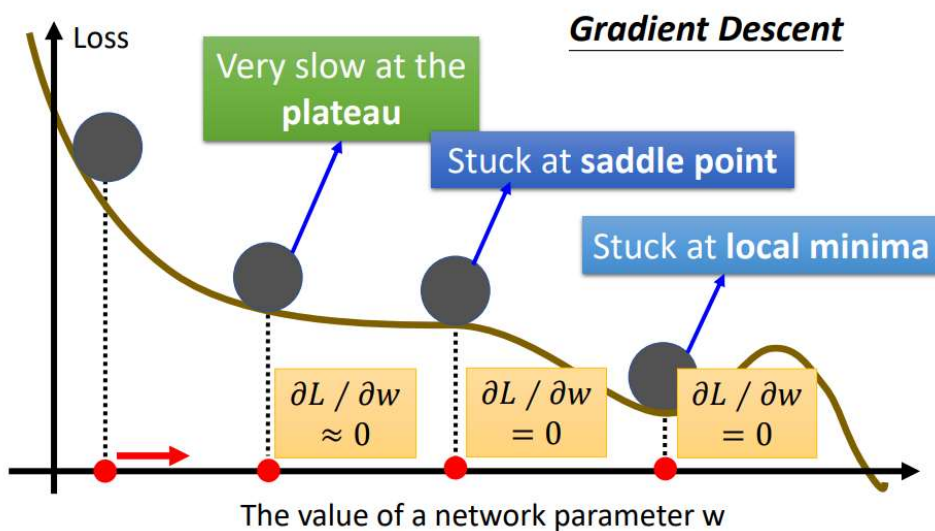
Empirical Study



$$\text{Minimum ratio} = \frac{\text{Number of Positive Eigen values}}{\text{Number of Eigen values}}$$

- 每一個點都代表訓練 network 訓練完之後，把它的 Hessian 拿出來進行計算，訓練到 gradient 很小，卡在 critical point，把那組參數拿出來分析，看它比較像是 saddle point 還是比較像 local minima
- 縱軸代表 training 時的 loss，loss 沒辦法再下降時
- 橫軸是 minimum ratio
- 如果所有的 eigen value 都是正的，代表 critical point 是 local minima
- 如果有正有負代表是 saddle point
- local minima 不常見

Small Gradient ...



Gradient 非常小的時候，有甚麼可能的解決辦法

課程網頁 (<https://speech.ee.ntu.edu.tw/~hylee/ml/2022-spring.php>).

tags: 2022 李宏毅_機器學習