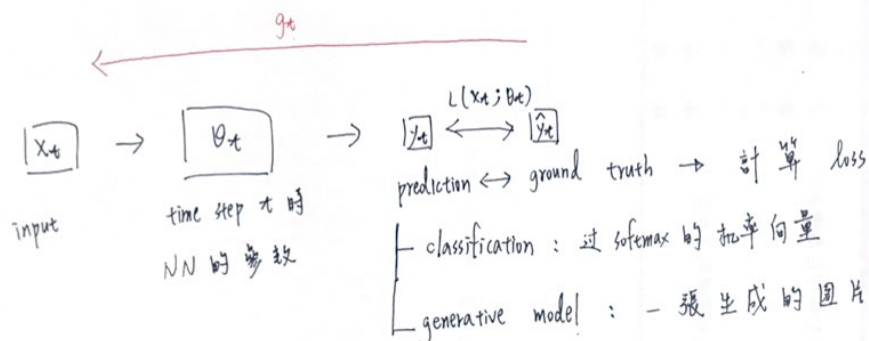


$\theta_t$  = model parameters  $\rightarrow$  此訓練想 optimize 的問題.

$\nabla L(\theta_t) / g_t$  : 計算某組  $D$  在特定時候 loss 的 gradient  $\rightarrow$  依  $\theta_t$  optimize 得  $\theta_{t+1}$  的 constant

$m_{t+1}$  : momentum, 記錄之前的 gradient  $\rightarrow \theta_t \rightarrow \theta_{t+1}$  的 gradient (會壓縮, 不會全概括)



### Optimization

找 -  $y$  使 training data 的  $x$  算出的 loss 最小  
 $\rightarrow$  parameters 越貼近訓練資料越好, 使  $y_t \approx \hat{y}_t$

On-line : 每  $t$  time step 只有 -  $y$   $x_t$

Off-line : 每  $t$  time step 可得

所有的  $x, y$ , ground truth  
容易  $\rightarrow$  get all training data

※ 無這麼大的 memory

$\rightarrow$  一次只碼 -  $y$   $x$

SGD : 於初始的  $\theta_0$  開始, 計算 gradient 使往反方向移動

→  $\therefore$  gradient 為  $L$  增加的方向,  $\therefore$  往反方向才可 find 最小  $L$  的  $\theta$

$$\theta_t = \theta_{t-1} - \eta \cdot g_{t-1}$$

SGD with Momentum : Movement  $v^0 = 0$

(SGDM) update → Movement.  $v^1 = \lambda v^0 - \eta \nabla L(\theta^0)$  → 如 慣性,  $\lambda$  中 慣性 影響力 成正比

SGD 算出 速率 update 的方向.

$v_i$  → 做時間的 weighted sum, 過去所有 gradient 的總和.

Adagrad : 防止一開始 gradient 過大, 導致移動過多而走向更差的位置.

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{\sum_{i=0}^{t-1} (g_i)^2}} g_{t-1}$$

time step gradient 的和, gradient 大 → 移動小, 避免跳過最小  
小 → 可放心移動

gradient 持續加總, 若一開始 gradient 太大, learning rate 過小, 可能走沒幾步便卡住.

RMSProp : 確保  $v_t$  不斷變化, gradient 避免持續變大, 使 Optimizer 不會卡住

(以 Momentum)

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{v_t}} g_{t-1}$$

$$v_1 = g_0^2$$

$$v_t = \alpha v_{t-1} + (1-\alpha)(g_{t-1})^2$$

$v_i$  → sum of gradient

無法處理 gradient = 0, 導致 Optimizer 卡住問題.

Adam : SGD  $\rightarrow \theta_t = \theta_{t-1} - \eta m_t$

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

+

RMSProp  $\rightarrow \theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{v_t}} g_{t-1}$

$$v_t = g_{t-1}^2$$

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

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

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \quad \beta_1 = 0.9$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t} \quad \beta_2 = 0.999$$

$$\epsilon = 10^{-8}$$

$\hat{m}_t \rightarrow$  確保  $m_t$  不會隨時間增加越來越大

$\therefore \beta_1 < 1 \quad \therefore m$  - 開始較小, 隨時間穩定

$\epsilon \rightarrow$  避免  $t=0, v_t=0$ , 造成  $\eta$  無限大  
epsilon