

* 建模基本流程

$$1. \hat{y} = f(w, x) \quad \rightarrow \theta$$

$$2. C = C(\hat{y}, y) \quad \rightarrow \text{MLE view: cross entropy, MSE}$$

$$3. w^* = \underset{w}{\operatorname{argmin}} C \quad \rightarrow \text{Difficult when it's deep}$$

↓
optimization method

$$\downarrow$$

$$x \rightarrow f(w^*, x) \xrightarrow{\text{predict}} \hat{y}$$

$$\downarrow$$

$$\text{Gradient descent}$$

↓

$$w_{t+1} = w_t - LR \nabla_w C$$

* 梯度下降法

國中：

$$\text{斜率} = \frac{\Delta y}{\Delta x}$$

<補充>

2D 示意圖

高中：

$$\text{一階導數} = \lim_{\Delta x \rightarrow 0} \frac{f(x+\Delta x) - f(x)}{\Delta x}$$

現在：

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \quad x \in \mathbb{R}^n$$

$$w_{t+1} = w_t - LR \nabla_w C ?$$

$$\phi = \lim_{h \rightarrow 0} \frac{C(w_t + h \cdot u) - C(w_t)}{h}$$

$$u = \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix} = u_1 e_1 + u_2 e_2 + \dots$$

by Taylor^{1st} at w_t

$$\|u\| = 1$$

$$C(w_t + hu) \approx C(w_t) + \left. \frac{\partial C}{\partial w_{t1}} \right|_{w=w_t} \cdot (w_{t1} + hu_1 - w_{t1}) + \dots$$

$$+ \left. \frac{\partial C}{\partial w_{tn}} \right|_{w=w_t} \cdot (w_{tn} + hu_n - w_{tn})$$

$$= C(w_t) + h \cdot \nabla C(w_t)^T \cdot u$$

代 $\lambda \phi$

$$\lim_{h \rightarrow 0} \frac{C(w_t + h \cdot u) - C(w_t)}{h} = \lim_{h \rightarrow 0} \frac{h \cdot \nabla C(w_t)^T u}{h} = \nabla C(w_t)^T u$$

$$\Rightarrow \phi = \|\nabla C(w_t)\| \cdot \cos \theta$$

$$\Rightarrow u^* = \arg \min_u \phi = \phi (\theta = \pi)$$

ex: Regression

$$1. \hat{y} = XW + b$$

$$W = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}_{d \times 1} \quad X = \begin{bmatrix} 1 & x_1 & \dots & x_d \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & \dots & x_d \end{bmatrix}_{N \times d}$$

$$2. C = \|\hat{y} - y\|^2$$

$$3. \nabla_W C, \nabla_b C = ?$$

$$b = \begin{bmatrix} b_1 \\ \vdots \\ b_N \end{bmatrix}_{N \times 1} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$$

$$C = (\hat{y} - y)^T (\hat{y} - y)$$

$$= (XW + b - y)^T \cdot (XW + b - y)$$

$$= ((XW)^T XW + (XW)^T b - (XW)^T y + b^T XW + b^T b - b^T y - y^T XW - y^T b + y^T y)$$

$$\nabla_W C = (2X^T XW + X^T b - X^T y + X^T b - X^T y)$$

$$= 2X^T (XW + b - y) = 2X^T (\hat{y} - y)$$

<補充>

1. Normal Equation?

2. 幾何意義?

$$\nabla_b^C = z(\hat{y} - y)$$

方法 2 :

$$\begin{aligned} \frac{\partial C}{\partial w_i} &= \sum_j^N \frac{\partial C}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial w_i} \\ &= \sum_j^N z(\hat{y}_j - y_j) \cdot x_{ji} \end{aligned}$$

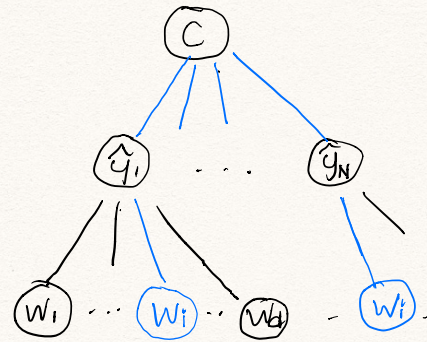
$$= z \cdot (\hat{y} - y) \odot x$$

\Downarrow

$$\nabla_w^C = z x^T (\hat{y} - y)$$

$$\nabla_b^C = z(\hat{y} - y)$$

Chain Rule



$$\hat{y}_j = \sum_k^d w_k x_{jk} + b_j$$

* 常見最佳化方法

(1) SGD:

$$w_{t+1} = w_t - LR \cdot \nabla_w C$$

<補充>

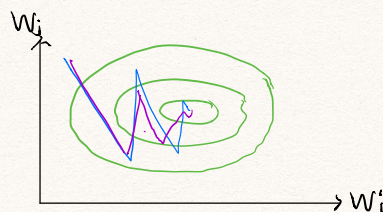
weight decay

(2) Momentum:

$$v_{t+1} = \rho v_t + G_{t+1}$$

$$w_{t+1} = w_t - LR \cdot (\rho v_t + G_{t+1})$$

$$v_0 = 0$$



(3) AdaGrad :

$$W_{t+1} = W_t - LR \cdot G_{t+1} \odot \frac{1}{\sqrt{r_{t+1}} + \epsilon}$$

$$r_{t+1} = \sum_z^{t+1} G_z \odot G_z$$

↙
elementwise operation

- LR 逐漸調小, 某方向太快讓它緩緩
- ↳ 有時是缺美, 收斂太快

(4) RMSprop :

- AdaGrad + Momentum Concept on r

$$W_{t+1} = W_t - LR \cdot G_{t+1} \odot \frac{1}{\sqrt{r_{t+1}} + \epsilon}$$

$$r_{t+1} = \underbrace{\rho}_{0.99 \text{ or } 0.9} r_t + (1 - \rho) G_{t+1} \odot G_{t+1}$$

(5) Adam :

- First moment

$$\hat{v}_{t+1} = \frac{\rho_1 v_t + (1 - \rho_1) G_{t+1}}{(1 - \rho_1^{t+1})}$$

<補充>

為什麼除以 $(1 - \rho_1^{t+1})$?

$$v_{t+1} = \rho_1 (\rho_1 v_{t+1} + \dots) + (1 - \rho_1) G_{t+1}$$

$$= (1 - \rho_1) \sum_z G_z \rho_1^{t-z}$$

$$E(v_{t+1}) = (1 - \rho_1^{t+1}) E(G_t)$$

$$E(G_t) = \frac{E(v_{t+1})}{1 - \rho_1^{t+1}}$$

- Second moment

$$\hat{r}_{t+1} = \frac{\rho_2 r_t + (1 - \rho_2) G_{t+1} \odot G_{t+1}}{(1 - \rho_2^{t+1})}$$

$$\Rightarrow W_{t+1} = W_t - LR \cdot \frac{\hat{v}_{t+1}}{\sqrt{\hat{r}_{t+1}} + \epsilon}$$