

梯度下降法

* 建模基本流程

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

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

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

↓
optimization method

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

↓
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：

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

現在：

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

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

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

by Taylor^{1st} at w_t

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

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

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

def $\lambda \phi$

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

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

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

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

$$\|u\| = 1$$

§ Regression 為例

$$1. \hat{y} = x w + b$$

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}_{d \times 1} \quad x = \begin{bmatrix} 1 & 2 & \dots & d \\ \hline \vdots & \vdots & \ddots & \vdots \\ \hline 1 & 2 & \dots & d \end{bmatrix}_{N \times d}$$

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

$$3. \nabla_w^C, \nabla_b^C = ?$$

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

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

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

$$= ((xw)^T xw + (xw)^T b - (xw)^T y + b^T xw + b^T b - b^T y \\ - y^T xw - y^T b + y^T y)$$

$$\nabla_w^C = (2x^T xw + x^T b - x^T y + x^T b - x^T y) \\ = 2x^T (xw + b - y) = 2x^T (\hat{y} - y)$$

<補充>

1. Normal Equation?

2. 幾何意義?

$$\nabla_b^C = 2(\hat{y} - y)$$

方法 2：

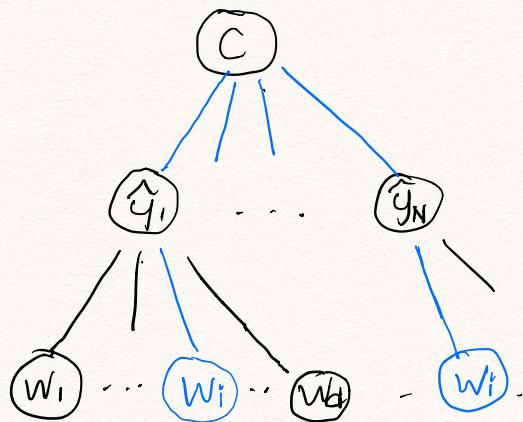
$$\begin{aligned}\frac{\partial C}{\partial w_i} &= \sum_j^N \frac{\partial C}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial w_i} \\ &= \sum_j^N 2(\hat{y}_j - y_j) \cdot x_{ji} \\ &= 2 \cdot (\hat{y} - y) \odot x\end{aligned}$$

¶

$$\nabla_w C = 2X^T(\hat{y} - y)$$

$$\nabla_b C = 2(\hat{y} - y)$$

Chain Rule



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

* $C = \|\hat{y} - y\|^2$ 怎麼來的？

$$y^{(i)} = X^{(i)}w + \underbrace{\varepsilon^{(i)}}_{\text{常態}}$$

$$P(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}G} \exp\left(-\frac{(E^{(i)})^2}{2G^2}\right)$$

$$P(y^{(i)} | X^{(i)}; w) = \frac{1}{\sqrt{2\pi}G} \exp\left(-\frac{(y^{(i)} - X^{(i)}w)^2}{2G^2}\right)$$

$$L(w) = \prod_{i=1}^N P(y^{(i)} | x^{(i)}; w)$$

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{(y^{(i)} - x^{(i)}w)^2}{2\sigma^2}\right)$$

$$\ell(w) = \log L(w)$$

$$= \log \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{(y^{(i)} - x^{(i)}w)^2}{2\sigma^2}\right)$$

$$= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{(y^{(i)} - x^{(i)}w)^2}{2\sigma^2}\right)$$

$$= N \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^N (y^{(i)} - x^{(i)}w)^2$$

$$C := \frac{1}{2} \sum_{i=1}^N (y^{(i)} - x^{(i)}w)^2$$

$$= \frac{1}{2} \|\hat{y} - y\|^2$$

* 常見最佳化方法

(1) SGD:

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

<補充>

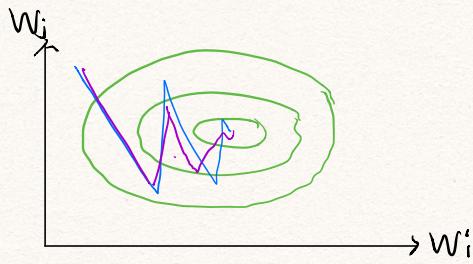
weight decay

(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:

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

$$(r_{t+1} = \sum_i^{} G_i \odot G_i)$$

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

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

$$(r_{t+1} = \rho r_t + (1-\rho) G_{t+1} \odot G_{t+1})$$

$\downarrow 0.99 \text{ or } 0.9$

(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 $(1 - \rho_1^{t+1})$?

- Second moment

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

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

$$= (1 - \rho_1) \sum G_t \rho_1^{t-1}$$

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

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

(6) Second-Order Optimization

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^T \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^T H(\theta - \theta_0)$$

$$\nabla_{\theta} J(\theta_0) + H(\theta - \theta_0) = 0$$

$$\Rightarrow \theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$