$$\ell_{DPO}(\phi) = -\mathbb{E}_{(x, y_w, y_t) \sim D}\left[\log \sigma(\delta(\phi))\right] \text{ where}$$

$$\text{where } \delta(\phi) = \beta \log \frac{\pi_\phi(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_\phi(y_t|x)}{\pi_{ref}(y_t|x)}$$

$$\rightarrow \ell_{DPO}(\phi) = -\mathbb{E}_{(x, y_w, y_t) \sim D}\left[\log \sigma(\delta(\phi))\right] \quad \text{——①}$$

$$\text{Gradient on both sides } \quad \nabla_\phi \ell_{DPO}(\phi) = \nabla_\phi \left(-\mathbb{E}_{(x, y_w, y_t) \sim D}\left[\log \sigma(\delta(\phi))\right]\right)$$

Properties of sigmoid $\rightarrow \sigma(z) = \dfrac{1}{1 + e^{-z}} \rightarrow \sigma'(z) = \sigma(z)(1 - \sigma(z))$

Derivative of $\log(x) = \dfrac{1}{x}$

$$\nabla_\phi \ell_{DPO}(\phi) = -\mathbb{E}_{(x, y_w, y_t) \sim D}\left[\frac{1}{\sigma(\delta(\phi))} \cdot \sigma(\delta(\phi))\left[1 - \sigma(\delta(\phi))\right] \cdot \nabla_\phi \delta(\phi)\right]$$

$$= -\mathbb{E}_{(x, y_w, y_t) \sim D}\left[\left[1 - \sigma(\delta(\phi))\right] \cdot \nabla_\phi \delta(\phi)\right]$$

$$\left(\text{Property of sigmoid} = \sigma(-\delta(\phi)) = 1 - \sigma(\delta(\phi))\right)$$

$$\nabla_\phi \ell_{DPO}(\phi) = -\mathbb{E}_{(x, y_w, y_t) \sim D}\left[\sigma(-\delta(\phi)) \cdot \nabla_\phi \delta(\phi)\right] \quad \text{——②}$$

ⓐ

Solving ⓐ $\quad \nabla_\phi \delta(\phi) = \nabla_\phi\left[\beta \log \frac{\pi_\phi(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_\phi(y_t|w)}{\pi_{ref}(y_t|w)}\right]$

$$= \nabla_\phi\left[\beta \log \pi_\phi(y_w|x) - \beta \log \pi_{ref}(y_w|x) - \beta \log \pi_\phi(y_t|x) + \beta \log \pi_{ref}(y_t|x)\right]$$

$$\nabla_\phi \delta(\phi) = \beta \frac{\nabla_\phi \pi_\phi(y_w|x)}{\pi_\phi(y_w|x)} - \frac{\beta \nabla_\phi \pi_\phi(y_t|x)}{\pi_\phi(y_t|x)} \quad \text{——ⓑ}$$

Putting ⓑ back in ②

$$\nabla_\phi \ell_{DPO}(\phi) = -\mathbb{E}_{(x, y_w, y_t) \sim D}\left[\sigma(-\delta(\phi))\left[\beta \frac{\nabla_\phi \pi_\phi(y_w|x)}{\pi_\phi(y_w|x)} - \frac{\beta \nabla_\phi \pi_\phi(y_t|x)}{\pi_\phi(y_t|x)}\right]\right]$$

Hence. Proved

$$\boxed{\nabla_\phi \ell_{DPO}(\phi) = \mathbb{E}_{(x, y_w, y_t) \sim D}\left[\sigma(-\delta(\phi))\left[-\beta \frac{\nabla_\phi \pi_\phi(y_w|x)}{\pi_\phi(y_w|x)} + \beta \frac{\nabla_\phi \pi_\phi(y_t|x)}{\pi_\phi(y_t|x)}\right]\right]}$$