

$$\nabla_{w_1}(J)$$

$$\nabla_{b_1}(J)$$

$$\nabla_{w_3} J = \nabla_{w_3} (L(\hat{y}, y) + \lambda \Omega(\theta))$$

$$= \nabla_{a_3} (L(\hat{y}, y)) \nabla_{a_3} (a_3) + \lambda \nabla_{w_3} (\Omega(\theta))$$

$$= g(\hat{y}) f'(a_3) h_2^T + \lambda \nabla_{w_3} (\Omega(\theta))$$

$$\nabla_{b_3} J =$$

$$g(\hat{y}) f'(a_3) + \lambda \nabla_{b_3} (\Omega(\theta))$$

Don't worry

I moved the text to the next page

$$\begin{aligned}\nabla_{h_2} [L(\hat{y}, y)] &= \nabla_{h_2} a_3 \nabla_{a_3} L(\hat{y}, y) \\ &= \nabla_{h_2} (\omega_3^T h_2 + b_3) (g \odot f'(a_3)) \\ &= \omega_3^T (g \odot f'(a_3))\end{aligned}$$

$$g \leftarrow \nabla_{h_2} (L(\hat{y}, y)) = \omega_3^T (g \odot f'(a_3))$$

$$\begin{aligned}\nabla_{a_2} (L(\hat{y}, y)) &= \nabla_{h_2} (L(\hat{y}, y)) \nabla_{a_2} h \\ &= g \nabla_{a_2} (f(a_2)) = g \cdot f'(a_2)\end{aligned}$$

$$\begin{aligned}\nabla_{\omega_2} \mathcal{J} &= \nabla_{\omega_2} (L(\hat{y}, y)) + \lambda \omega_2 (-\Omega(\theta)) \quad (2_{\text{layer}}) \\ &= \nabla_{a_2} (L(\hat{y}, y)) \nabla_{\omega_2} (a_2) + \lambda \nabla_{\omega_2} (-\Omega(\theta)) \\ &= (g \cdot f'(a_2)) (\nabla_{\omega_2} (\omega_2^T h_1 + b_2)) + \lambda \nabla_{\omega_2} (-\Omega(\theta))\end{aligned}$$

$$\begin{aligned}\nabla_{b_2} (\mathcal{J}) &= \nabla_{b_2} (L(\hat{y}, y) + \lambda \Omega(\theta)) \\ &= \nabla_{a_2} (L(\hat{y}, y)) \nabla_{b_2} (a_2) + \nabla_{b_2} (\lambda \Omega(\theta))\end{aligned}$$

$$= (g \cdot f'(a_2)) \nabla_{b_2} (\omega_2^T h_1 + b_2) + \nabla_{b_2} (\lambda \Omega(\theta))$$

$$= (g \cdot f'(a_2)) + \lambda \nabla_{b_2} (-\Omega(\theta))$$

$$\nabla_{h_1} (L(\hat{y}, y)) = \nabla_{a_2} (L(\hat{y}, y)) \nabla_{h_1} (a_2)$$

$$= (g \cdot f'(a_2)) \nabla_{h_1} (\omega_2^T h_1 + b_1)$$

$$= (g \cdot f'(a_2)) \omega_2^T$$

$$g \leftarrow \nabla_{h_1} (L(\hat{y}, y))$$

111⁴ $\nabla_{w_1}(\mathcal{J})$ can be found
 $\nabla_{b_1}(\mathcal{J})$ thru chain rule

in textbook \leftarrow he takes gradient
of loss function with regularization

$g \leftarrow \nabla_{\hat{y}}(\mathcal{J})$ this is incorrect!

$$\begin{aligned}\nabla_{\hat{y}}(\mathcal{J}) &= \nabla_{\hat{y}}(L(y, \hat{y}) + \lambda \Omega(w, b)) \\ &= \nabla_{\hat{y}}(L(y, \hat{y})) + \lambda \nabla_{\hat{y}}(\Omega(w, b))\end{aligned}$$

both $\hat{y} \leftarrow$
& Ω are functions of w_3 !!

This gradient is not zero!

HOWEVER

You may say that TB considers
regularization as part of the
loss function i.e. $L(\hat{y}, y)$ has
the regularized term inside

eg: $L(\hat{y}, y) = \text{MSE} + \lambda \Omega(w)$