

# Task 1

$$a) \quad C(\omega) = \frac{1}{N} \sum_{n=1}^N C^n$$

$$C^n(\omega) = -(\hat{y}^n \ln(\hat{y}^n) + (1 - \hat{y}^n) \ln(1 - \hat{y}^n))$$

$$\frac{\partial f(x^n)}{\partial w_i} = x_i^n f(x^n) (1 - f(x^n))$$

$$\frac{\partial C^n}{\partial w_i}$$

$$\frac{\partial C^n(\omega)}{\partial w_i} = x_i^n C^n(\omega) (1 - C^n(\omega))$$

$$b) \frac{\partial \hat{C}(\omega)}{\partial \omega_{kj}} = \frac{\partial C}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w}$$

$$\frac{\partial \mathcal{L}}{\partial a} \bigg|_{a = \frac{1}{f(x)}} = \frac{\partial \left( - \sum_{k=1}^K y_k^n \ln(\hat{y}_k^n) \right)}{\partial \hat{y}_k^n}$$

$$i = 0, \quad \frac{\partial y_i}{\partial x_i} = \frac{e^{x_i} \sum_k e^{x_k} - e^{x_i} e^{x_i}}{(\sum_k e^{x_k})^2} = y_i (1 - y_i)$$

$$i \neq j, \quad \frac{\partial y_i}{\partial x_j} = \frac{c - e^{x_j} e^{x_i}}{(\sum_k e^{x_k})^2} = -\frac{1}{y_i} \frac{1}{y_j}$$

$$\frac{\partial C}{\partial y_i} = - \sum_j y_j \frac{1}{y_j}$$

$$\frac{\partial \mathcal{L}}{\partial x_j} = - \sum_{i: i \neq j} y_i \frac{1}{y_i} \frac{\partial y_i}{\partial x_j} - y_j \frac{1}{y_j} \frac{\partial y_j}{\partial x_j}$$

$$= - \sum_{i=1}^n \gamma_i \frac{1}{\gamma_i} \left( -\frac{1}{\gamma_i} (-\frac{1}{\gamma_i} \gamma_i) - \gamma_i \frac{1}{\gamma_i} \frac{1}{\gamma_i} (1 - \frac{1}{\gamma_i}) \right)$$

$$\Rightarrow \frac{\partial C}{\partial a} = -x_0 (y_1 - \hat{y}_1)$$

4a)

$$R(w) = \|w\|^2 = \frac{1}{2} \sum_{i,j} w_{i,j}^2$$

$$J = C(w) + \lambda R(w)$$

$$\frac{\partial J}{\partial w} = \frac{\partial C}{\partial w} + \underbrace{\frac{\partial \left( \lambda \frac{1}{2} \sum_{i,j} w_{i,j}^2 \right)}{\partial w}}_{\substack{\text{create fun!} \\ \nearrow}} = \frac{\partial C}{\partial w} + \lambda \cdot w$$

$X_{\text{train}}$  : images for training

$Y_{\text{train}}$  : labels for training

$\text{Image 1} = \text{image1.reshape}(28, 28)$

$\text{plt.imshow}(\text{Image 1})$

$\hat{y}$  : Prediction

$y$  : ground truth

Gradient descent:

$$W_{b+1} = W_b - \alpha \frac{\partial \hat{L}(w)}{\partial w}$$

loss plot

Y-axis : avg loss

X-axis : Number of training steps

$$X : 1005 \times 785 \rightarrow$$

$\nearrow$   $\nwarrow$   
 Feder Kegel

$$W : 785 \times 1 \rightarrow$$

$$W \times X = 785 \times 1 \times 1005 \times 785$$

$X$   $W$   $Feder$

$$W^T X = 1 \times 785 \times 1005 \times 785$$