

1.2 (b)

| Layer | Input | Output |
|----------|------------------------------|---|
| Linear 1 | x | $w^{(1)}x + b^{(1)} = z_1$ |
| f | $w^{(1)}x + b^{(1)}$ | $\max(0, w^{(1)}x + b^{(1)}) = \hat{y}_1 = z_2$ |
| Linear 2 | \hat{y}_1 | $w^{(2)}\hat{y}_1 + b^{(2)} = z_3$ |
| g | $w^{(2)}\hat{y}_1 + b^{(2)}$ | $w^{(2)}\hat{y}_1 + b^{(2)} = \hat{y}$ |
| Loss | \hat{y} | $\ \hat{y} - y\ _2^2$ |

(c)

$$\frac{dl}{dz_3} = \frac{dl}{d\hat{y}} \quad (\text{identity act.})$$

Assuming some dimension.

$$x \in \mathbb{R}^{d_{in}}$$

$$\frac{dl}{dw^{(2)}} = \frac{dl}{dz_3} \cdot \frac{dz_3}{dw^{(2)}} = \frac{dl}{d\hat{y}} \cdot (z_2)^T$$

$$z_2, z_1, b^{(1)} \in \mathbb{R}^{d_{2 \times 1}}$$

$$w^{(1)} \in \mathbb{R}^{d_2 \times d_1}$$

$$z_3, b^{(2)}, \hat{y} \in \mathbb{R}^{d_3 \times 1}$$

$$w^{(2)} \in \mathbb{R}^{d_3 \times d_2}$$

$$\frac{dl}{db^{(2)}} = \frac{dl}{dz_3} = \frac{dl}{d\hat{y}}$$

$$\text{Now, } \frac{dl}{dz_2} = \frac{dl}{dz_3} \cdot \frac{dz_3}{dz_2} = w_2^T \cdot \frac{dl}{d\hat{y}}$$

$$\frac{dl}{dz_1} = \frac{dl}{dz_2} \odot \frac{dz_2}{dz_1} = \left(w_2^T \cdot \frac{dl}{d\hat{y}} \right) \odot (z_1 > 0) \quad \leftarrow \text{hadamard product.}$$

$$\frac{dl}{dw^{(1)}} = \frac{dl}{dz_1} \cdot \frac{dz_1}{dw^{(1)}} = \frac{dl}{dz_1} \cdot x^T$$

$$\frac{dl}{db^{(1)}} = \frac{dl}{dz_1}$$

writing all here,

$$\frac{dl}{dw^{(2)}} = \frac{dl}{d\hat{y}} \cdot \frac{d\hat{y}}{dz_3} \cdot z_2^T = \frac{dl}{d\hat{y}} \cdot \frac{d\hat{y}}{dz_3} \cdot [(w^{(1)}x + b^{(1)})^+]^T$$

$$\frac{dl}{db^{(2)}} = \frac{dl}{d\hat{y}}$$

$$\frac{dl}{dw^{(1)}} = \left[w_2^T \frac{dl}{d\hat{y}} \odot (w^{(1)}x + b^{(1)} > 0) \right] \cdot x^T$$

$$\frac{dl}{db^{(1)}} = \left[w_2^T \frac{dl}{d\hat{y}} \odot (w^{(1)}x + b^{(1)} > 0) \right]$$

$$(d) \quad \frac{dl}{d\hat{y}} = \begin{bmatrix} dl/d\hat{y}_1 \\ \vdots \\ dl/d\hat{y}_{d_3} \end{bmatrix} \quad \frac{d\hat{y}}{dz_3} = I_{d_3 \times d_3} \quad (\text{diff identity})$$

$$\frac{dz_2}{dz_1} = \begin{matrix} \text{ReLU grad.} \\ (z_1 > 0) \end{matrix} = (w^{(1)}x + b^{(1)} > 0)_{d_2}$$

Ideally non differentiable

1.3. (a) with $f = g = \sigma$

→ the forward pass eqns change accordingly.

$$\rightarrow \sigma'(r) = \sigma(r) (1 - \sigma(r)) \quad \leftarrow \text{elementwise product}$$

↑ derivative changes

$$(b) \quad l(\hat{y}, y) = -\frac{1}{n} \sum_{i=1}^n y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)$$

$$\frac{dl}{d\hat{y}_i} = -\frac{1}{n} \left[\frac{y_i}{\hat{y}_i} - \frac{(1 - y_i)}{(1 - \hat{y}_i)} \right]$$

→ make vectors, just derivative changes.

(c) using a σ for intermediate activations

results in a high probability that the

gradients will explode or die.

