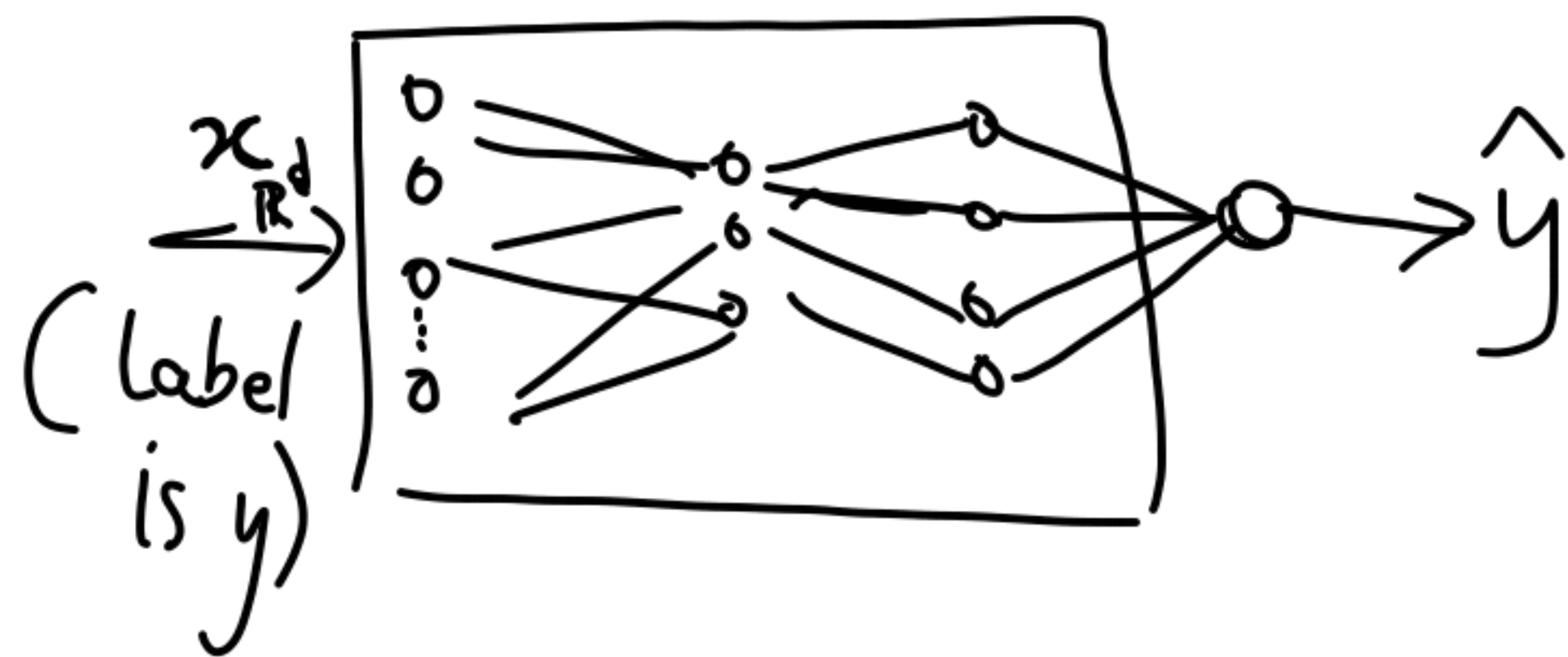


CS 725

Loss functions typically used in training neural networks.

① REGRESSION

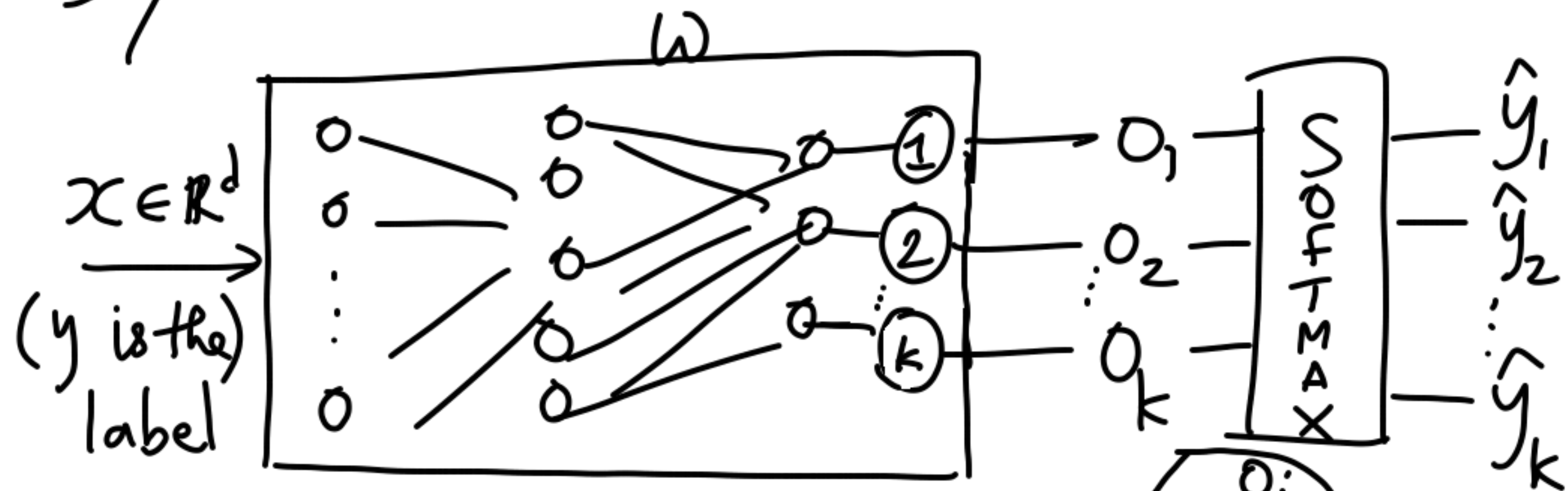


$$L(\hat{y}, y; w) = \frac{1}{2} (\hat{y} - y)^2$$

$\nearrow f_{NN}(x; w)$

[Convenience when computing
the gradient of L]

B) CLASSIFICATION



For a k-class
Classification
Task

$$\frac{e^{o_j}}{\sum_j e^{o_j}}$$

Predicted Probability Distribution

CROSS-ENTROPY
LOSS

$$L(\hat{y}, y, w) = \sum_{k=1}^K -y_k \log \hat{y}_k$$

$$[\hat{y}_1, \dots, \hat{y}_k]$$

y_k is either
0 or 1
 \hat{y}_k is a
probability

Training a neural network

Neural networks: Loss function is non-convex

Obvious choice to optimize the neural network loss function:

$\theta = \{w_0, w_1, \dots, w_N\}$ is set of all parameters of a neural network

GRADIENT DESCENT

$\theta \leftarrow$ initialize the weight vector

While (not stopping criterion) do

Pick a single example $(x, y) \in \mathcal{D}$ // SGD

$L \leftarrow L(f_\theta(x), y)$

end while

$\theta \leftarrow \theta - \eta \nabla L$

CRITICAL ASPECTS OF

NN TRAINING:

- ① Wt initialization
- ② Stopping criteria
- ③ Learning rate

For (S)GD training, we need to compute $\nabla_{\theta} L$

which is $\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_N}$

This gradient is computed efficiently using the BACKPROPAGATION
ALGORITHM

Preliminaries

Univariate chain rule ; Given a function $f(y)$, $y = g(x)$

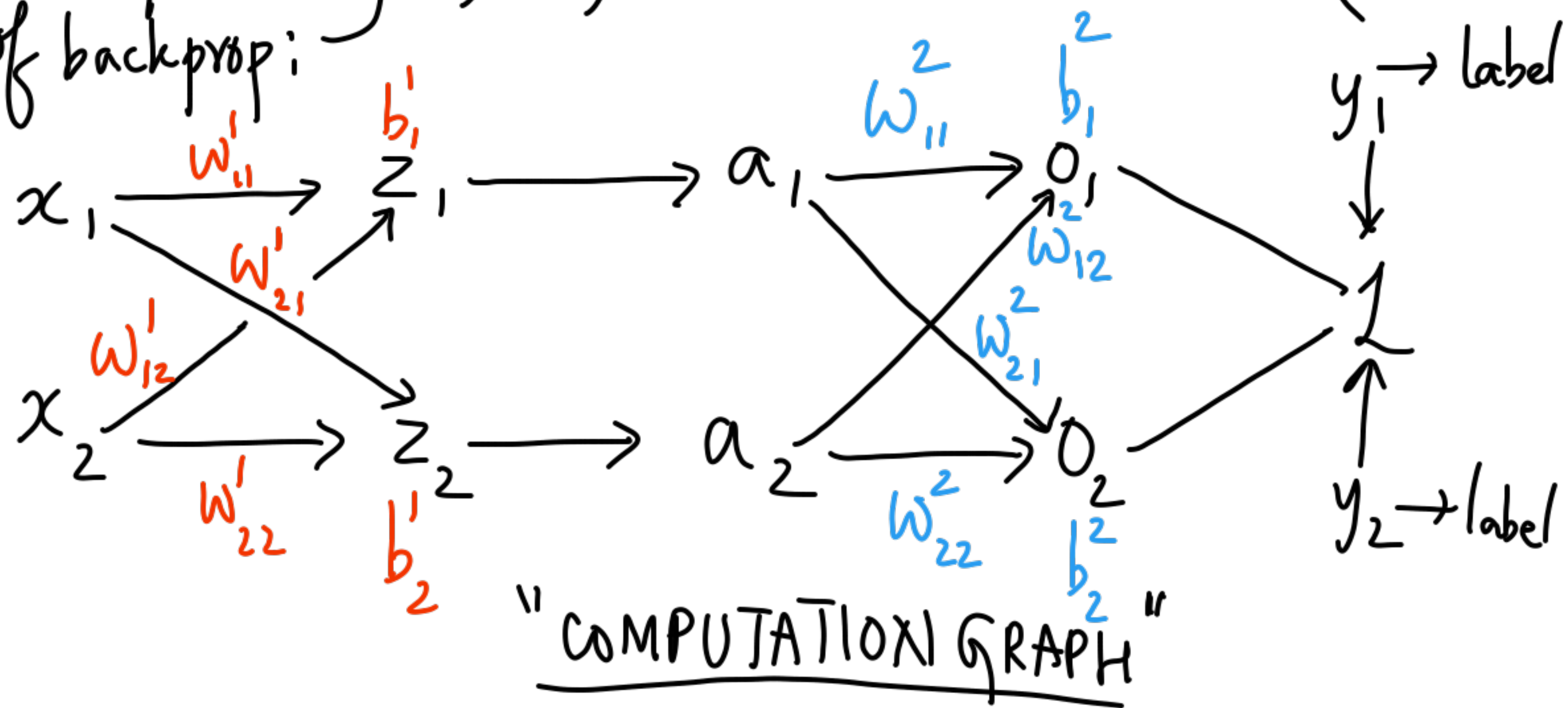
$$\frac{df}{dx} = \frac{df}{dy} \cdot \frac{dy}{dx}$$

Multivariate chain rule ; Given $f(y_1, \dots, y_n)$, $y_i = g_i(x)$

$$\frac{\partial f}{\partial x} = \sum_i \frac{\partial f}{\partial y_i} \cdot \frac{\partial y_i}{\partial x}$$

Backpropagation: Two-pass algorithm

Example using a feedforward neural network (multi-layer perceptron, MLP) of backprop:



① Forward pass

$$z_i = \sum_j w'_{ij} x_j + b'_i$$

$$a_i = \sigma(z_i)$$

$$o_k = \sum_i w^2_{ki} a_i + b^2_k$$

$$L = \frac{1}{2} \sum_k (o_k - y_k)^2$$

② BACKWARD PASS

$$\frac{\partial L}{\partial L} = 1 \quad (\text{Base case})$$

$$\frac{\partial L}{\partial o_k} = \frac{\partial L}{\partial L} \cdot \frac{\partial L}{\partial o_k} = o_k - y_k$$

$$\frac{\partial L}{\partial w_{ki}^2} = \frac{\partial L}{\partial o_k} \cdot \frac{\partial o_k}{\partial w_{ki}^2} = \frac{\partial L}{\partial o_k} \cdot a_i$$

$$\frac{\partial L}{\partial b_k^2} = \frac{\partial L}{\partial o_k} \cdot \frac{\partial o_k}{\partial b_k^2} = \frac{\partial L}{\partial o_k}$$

$$\frac{\partial L}{\partial a_i} = \sum_k \frac{\partial L}{\partial o_k} \cdot \frac{\partial o_k}{\partial a_i} = \sum_k \frac{\partial L}{\partial o_k} \cdot w_{ki}^2$$

$$\frac{\partial L}{\partial z_i} = \frac{\partial L}{\partial a_i} \cdot \frac{\partial a_i}{\partial z_i} = \frac{\partial L}{\partial a_i} \cdot \sigma'(z_i)$$

$$\frac{\partial L}{\partial w_{ij}^1} = \frac{\partial L}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}^1} = \frac{\partial L}{\partial z_i} \cdot x_j$$

$$\frac{\partial L}{\partial b_i^1} = \frac{\partial L}{\partial z_i} \cdot \frac{\partial z_i}{\partial b_i^1} = \frac{\partial L}{\partial z_i}$$