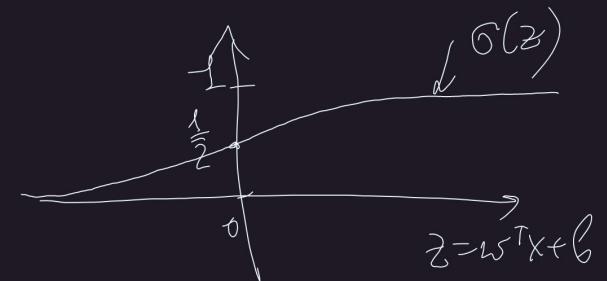


$$\left\{ (x_i, y_i) \right\}_{i=1}^N \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, +1\}$$

$y_i \in \{1, 2, \dots, K\}$



Decision rule:  $y(x) = \text{sign}(w^T x + b)$

$1 \times d \quad d \times 1 \quad 1 \times 1$

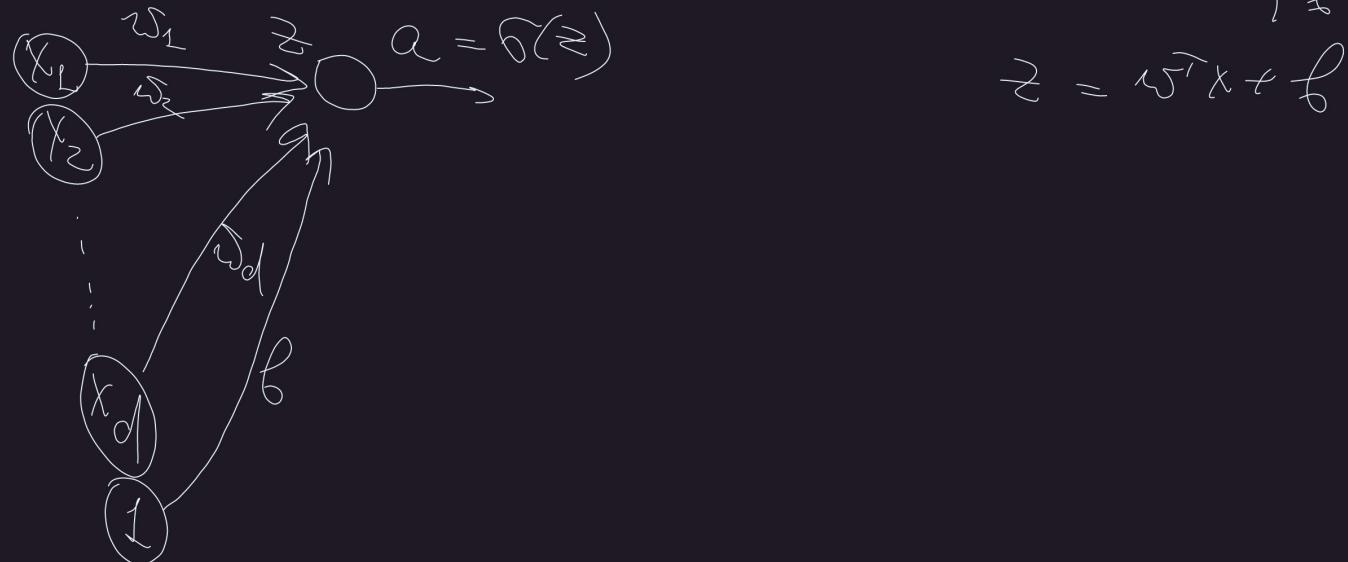
$$y(x) = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} (w_k^T x + b_k)$$

$$P(y = +1 \mid w, b, x) = \frac{1}{1 + \exp(-w^T x - b)} = \sigma(w^T x + b)$$

$$- P(y = k \mid w, b, x) = \frac{\exp(w_k^T x + b_k)}{\sum_{j=1}^K \exp(w_j^T x + b_j)} = \text{SoftMax}(w^T x + b)$$

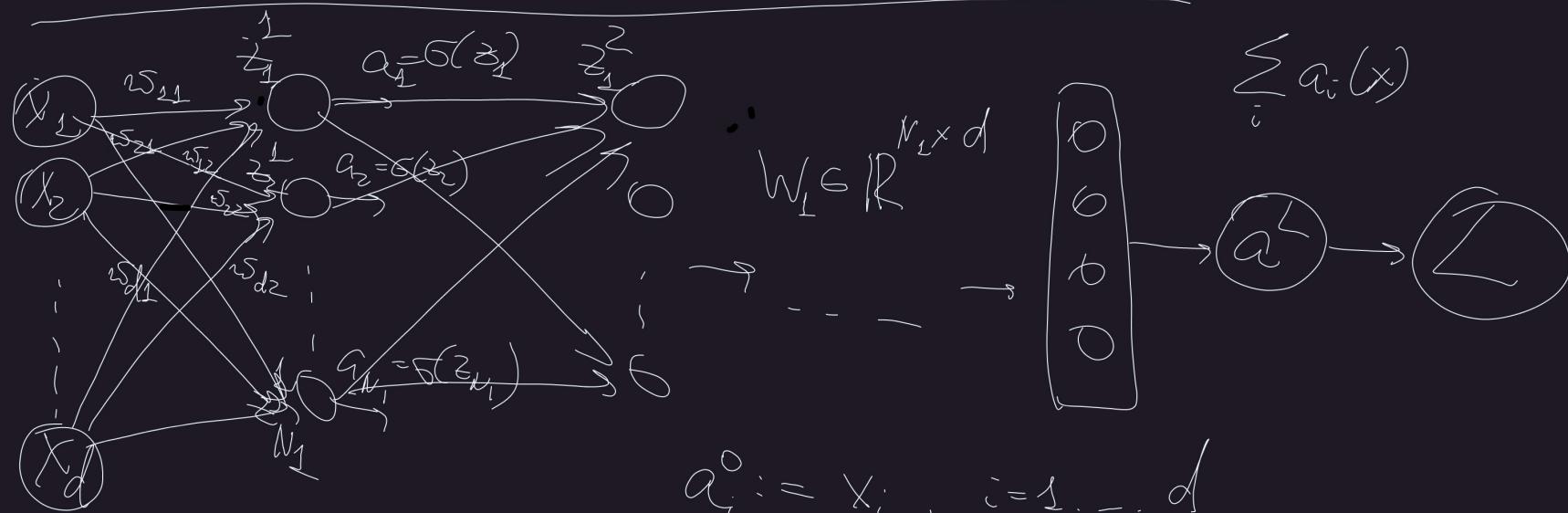
$$P(y_1, \dots, y_N | X, \omega, \theta) = \prod_{i=1}^N P(y_i | X_i, \omega, \theta) \rightarrow \max_{\omega, \theta}$$

$$\begin{aligned} -\log P(y_1 \sim y_N | X, \omega, \theta) &= -\sum_{i=1}^N \log P(y_i | X_i, \omega, \theta) \\ &= \sum_{i=1}^N \mathcal{L}(y_i, \omega^\top X_i + \theta) \rightarrow \min_{\omega, \theta} \end{aligned}$$



$$z_2^L = \omega_{12} \cdot x_1 + \omega_{22} \cdot x_2 + \dots + \omega_{d2} \cdot x_d + b_2$$

Multi-Layer Perceptron (MLP)

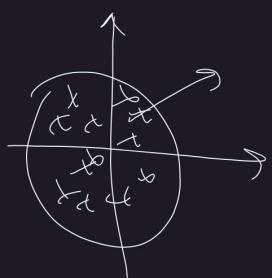
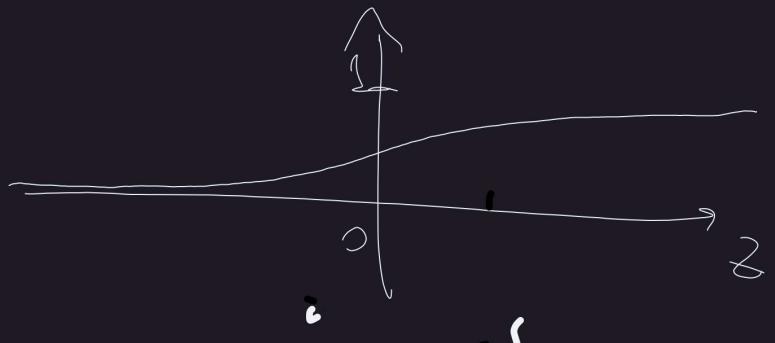


$$F(W) = \sum_{i=1}^N \left( y_i - a^L(x_i, W) \right)^2$$

$a_i^0 := x_i, \quad i = 1, \dots, d$   
 for  $\ell = 1, \dots, L$ :  
 $a^\ell = W^\ell a^{\ell-1} + b^\ell$   
 Output  $a^L = g(z^L)$

Activation functions:

$$\textcircled{L} \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$



$$\nabla_{\alpha^L} L = \left\{ \sum_{j=i+1}^L (W^{ij})^\top \frac{\partial a_j}{\partial z_j} \cdot \nabla_{\alpha^L} \right\}$$



$$\left\| \nabla_{\alpha^L} L \right\| \leq \prod_{j=i+1}^L \|W^{ij}\| \cdot \left\| \frac{\partial a_j}{\partial z_j} \right\| \cdot \left\| \nabla_{\alpha^L} \right\|$$

$$\text{diag}(\sigma'(z))$$

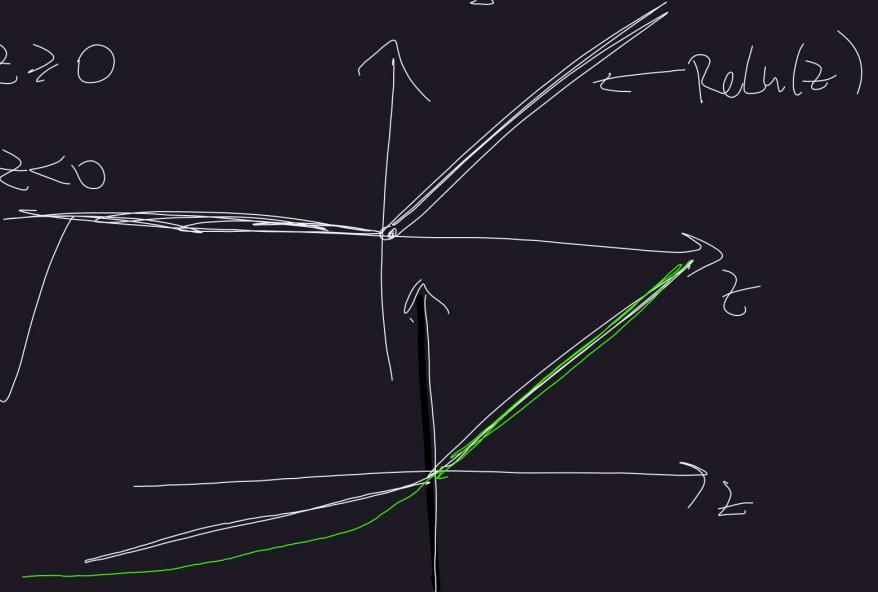
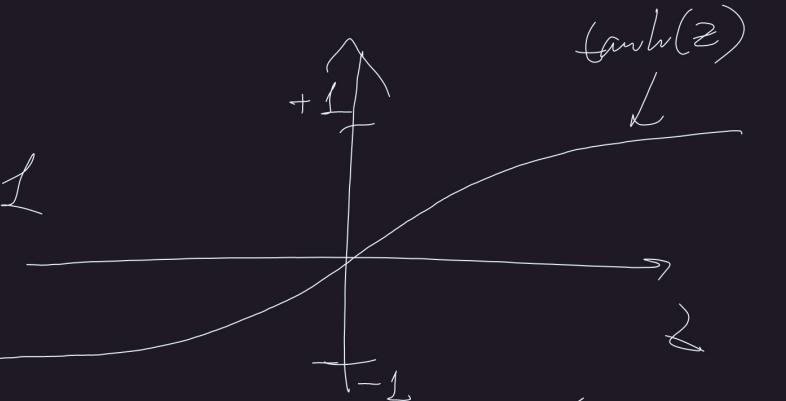
$$② g(z) = \tanh(z) = 2 \operatorname{csch}(2z) - 1$$

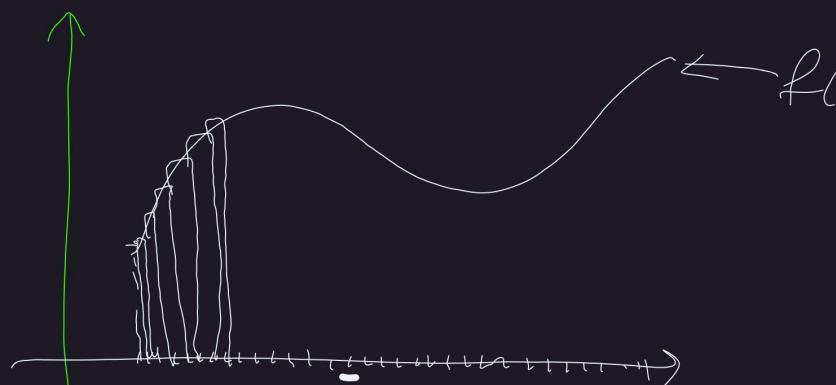
$$③ g(z) = \text{ReLU}(z) = \max(0, z)$$

$$④ g(z) = \text{Leaky ReLU}(z) = \begin{cases} z, & z \geq 0 \\ \alpha z, & z < 0 \end{cases}$$

$\alpha \in (0, 1)$

$$⑤ g(z) = E(u(z))$$



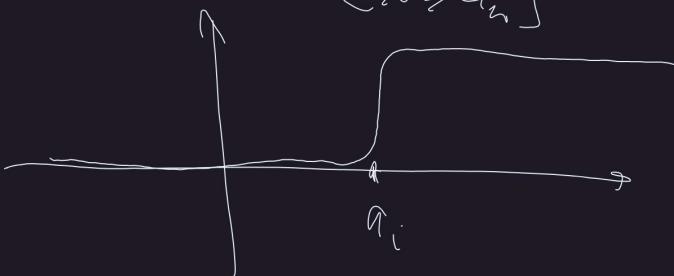


$$f(x) \approx \sum_i I[x \in [q_i, q_{i+1}]] \cdot f(q_i)$$

$$[x \geq q_i] \approx \delta(k(x - q_i))$$

$$[x \leq q_{i+1}] = 1 - [x \geq q_{i+1}]$$

$$\begin{aligned} I[q_i \leq x \leq q_{i+1}] &= \\ &= \left[ [x \geq q_i] + [x \leq q_{i+1}] \geq \frac{3}{2} \right] \\ &\quad \xrightarrow{x \geq q_i} \xrightarrow{x \leq q_{i+1}} \\ &\quad \xrightarrow{\vdots} \quad \xrightarrow{\vdots} \\ &\quad \xrightarrow{\oplus} \end{aligned}$$



# Automatic Differentiation

$$f(x)$$

$$\text{Find: } \nabla_x f(x)$$

6-th pos.

$$[0_{-}, 0, 1, 0_{-}, 0]$$

"

Finite Diff:

$$\frac{\partial f}{\partial x_i} \approx \frac{f(x + \varepsilon e_i) - f(x - \varepsilon e_i)}{2\varepsilon}$$

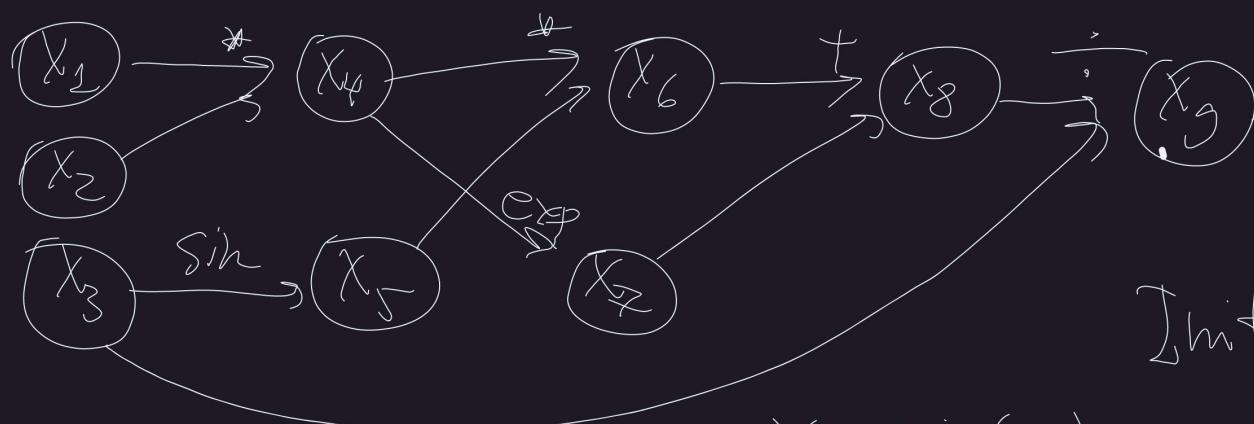
$$i = 1, \dots, d$$

$$\varepsilon = \varepsilon_m^{1/3}$$

$$f(x_1, x_2, x_3) =$$

$$\frac{\partial f}{\partial x_1} = \frac{\partial f}{\partial x_6} \cdot \frac{\partial x_6}{\partial x_1} + \frac{\partial f}{\partial x_7} \cdot \frac{\partial x_7}{\partial x_1}$$

$$x_1 \cdot x_2 \cdot \sin(x_3) + \exp(x_1 \cdot x_2)$$



Forward mode

$$\frac{\partial x_j}{\partial x_i}, \quad j=1 \dots 9, \quad i=1 \dots 3$$

$$x_5 = \sin(x_3)$$

$$\frac{\partial x_5}{\partial x_1} = \cos(x_3) \cdot \frac{\partial x_3}{\partial x_1}$$

$$x_4 = x_2 \cdot x_2$$

$$\frac{\partial x_4}{\partial x_1} = \frac{\partial x_2}{\partial x_1} \cdot x_2 + x_2 \cdot \frac{\partial x_2}{\partial x_1}$$

Init:

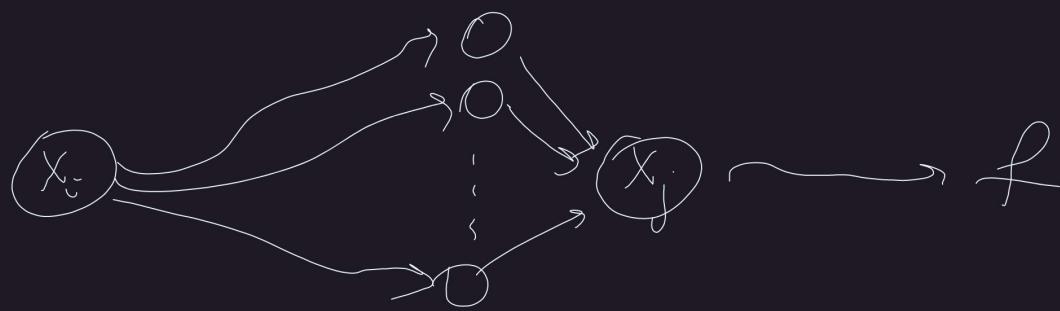
$$\frac{\partial x_L}{\partial x_1} = 1, \quad \frac{\partial x_L}{\partial x_2} = 0,$$

$$\frac{\partial x_L}{\partial x_3} = 0$$

$$\text{Init: } \frac{\partial f}{\partial x_g} = 1$$

$$\frac{\partial f}{\partial x_g} = \frac{\partial R}{\partial x_g} \cdot \frac{\partial x_g}{\partial x_g} = 1$$

$$= \left( \frac{\partial f}{\partial x_g} \right) \cdot \frac{1}{x_3}$$

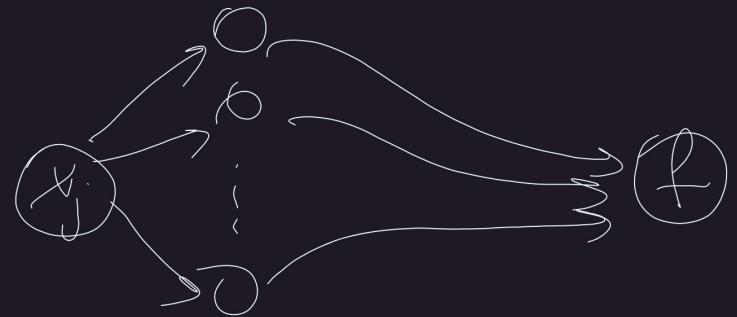


$$\frac{\partial \hat{f}}{\partial x_i} = \sum_{k: (k, i) \in \xi} \frac{\partial x_j}{\partial x_k} \left( \frac{\partial x_k}{\partial x_i} \right)$$

from prev. it.

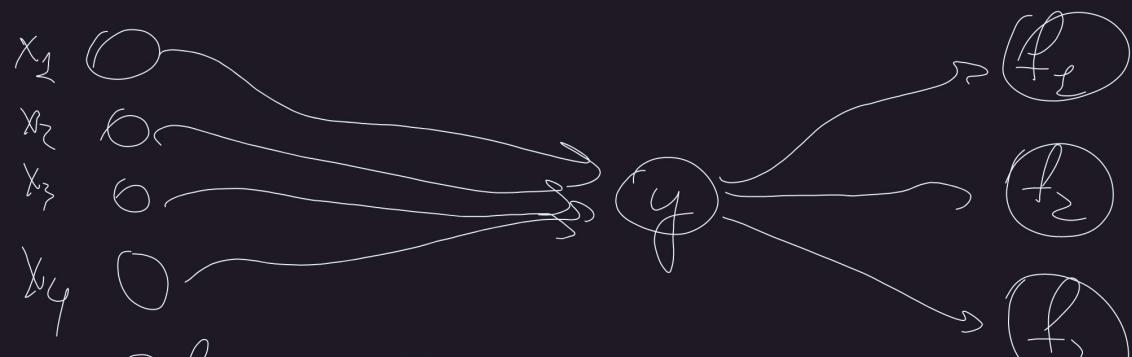
Backward mode

$$\frac{\partial \hat{f}}{\partial x_j} \quad j = g, p, \dots, 1$$



$$\frac{\partial \hat{f}}{\partial x_j} = \sum_{(j, k) \in \xi} \frac{\partial f}{\partial x_k} \cdot \frac{\partial x_k}{\partial x_j}$$

Cross Mode



$$\frac{\partial f_i}{\partial x_j} = \frac{\partial y}{\partial x_j} \cdot \frac{\partial f_i}{\partial y}$$

$$\left( \frac{\partial f}{\partial x} \right)_{ij} = \frac{\partial f_i}{\partial x_j}$$

$$f(x) = \text{tr}(Ax^T B) : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$$

$$\frac{\partial f}{\partial x} = \frac{\partial}{\partial x} \text{tr}(Ax^T B) = \text{tr}\left(\underbrace{\frac{\partial}{\partial x} Ax^T}_{n \times n} \underbrace{B}_{n \times n}\right)$$

Def  $f: U \rightarrow V$

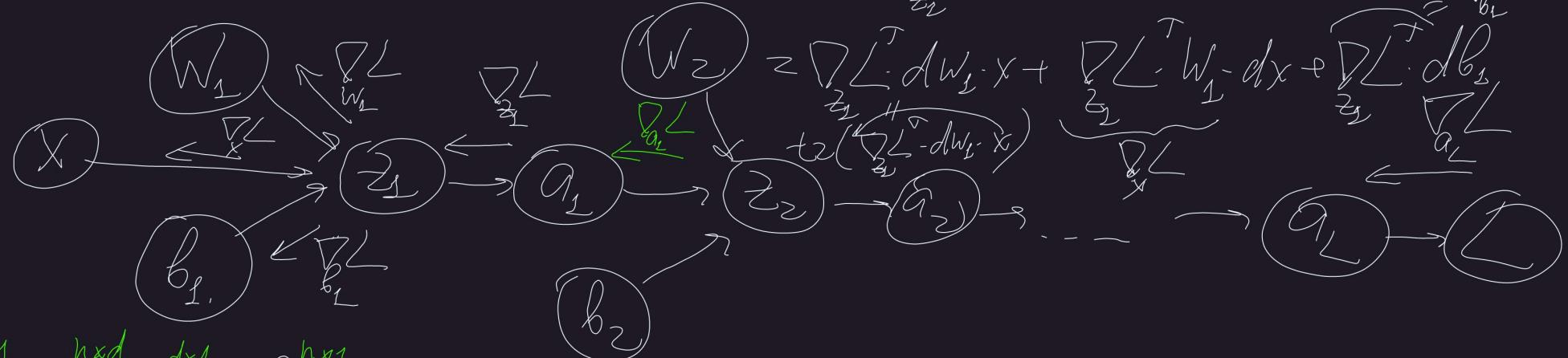
$$f(x+h) - f(x) = df(x)[h] + \tilde{O}(\|h\|)$$

$f: \mathbb{R} \rightarrow \mathbb{R}$	$\dot{df(x)}[h] = f'(x) \cdot h$	$\xrightarrow{\text{linear op. w.r.t. } h}$	$df(x)[h] =$
$f: \mathbb{R}^d \rightarrow \mathbb{R}$	$\dot{df(x)}[h] = \nabla f(x)^T h$		$f: (\mathbb{R}^{n \times n})^m \rightarrow \mathbb{R} = \langle \nabla_x f(x), h \rangle =$
			$= \text{tr}(\nabla_x f(x)^T h)$

$$f(x) = \text{tr}(A x^{-1} B)$$

$$\begin{aligned} df &= d\text{tr}(Ax^{-1}B) = \text{tr}(d(Ax^{-1}B)) = \\ &= \text{tr}(A \cdot \underbrace{dx^{-1}}_{-X^{-1}dX} \cdot B) = \text{tr}\left(-A \cdot X^{-1} dX \cdot \underbrace{X^{-1}B}_{-X^{-1}dX \cdot X^{-1}}\right) = \\ &\quad -X^{-1}dX \cdot X^{-1} \\ &\quad = \text{tr}(\nabla f(x)^T dX) = \end{aligned}$$

$$\nabla f(x) = -X^{-1} A^T B^T X^{-1} = \underbrace{\text{tr}(-X^{-1} B \cdot A \cdot X^{-1} dX)}_{\nabla f(x)^T}$$



$$z_1 = W_1 \cdot x + b_1 \quad \begin{matrix} n \times 1 \\ h \times d \\ d \times 1 \\ h \times 1 \end{matrix}$$

$$q_1 = g(z_1)$$

$$\begin{aligned} dz_1 &= d(W_1 x + b_1) = d(W_1 x) + db_1 = \\ &= dW_1 \cdot x + W_1 \cdot dx + db_1 \end{aligned}$$

$$\begin{aligned} dL &= \nabla_{z_1}^T (dW_1 \cdot x + W_1 \cdot dx + db_1) = \nabla_{b_1}^T \\ z_1 &= \nabla_{z_1}^T dL \cdot dW_1 \cdot x + \nabla_{W_1}^T W_1 \cdot dx + \nabla_{b_1}^T db_1 \\ &\quad + t_2(\nabla_{z_1}^T - dW_1 \cdot x) \quad \begin{matrix} n \times 1 \\ 1 \times d \\ n \times 1 \\ h \times 1 \end{matrix} \\ dL &= \nabla_{z_1}^T dL \cdot x^T \quad \begin{matrix} n \times 1 \\ h \times 1 \end{matrix} \\ \nabla_{W_1}^T &= W_1^T \nabla_{z_1}^T \quad \begin{matrix} n \times 1 \\ d \times 1 \end{matrix} \\ \nabla_{x}^T &= W_1^T \nabla_{z_1}^T \quad \begin{matrix} n \times 1 \\ d \times n \end{matrix} \end{aligned}$$

$$a_1 = g(z_1)$$

$$da_1 = \frac{\partial a_1}{\partial z_1} \cdot dz_1$$

$$\text{diag}(g'(z))$$

$$dL = \nabla_{a_1}^T \cdot da_1 = \underbrace{\nabla_{a_1}^T}_{\nabla_{z_1}} \cdot \underbrace{\frac{\partial a_1}{\partial z_1} \cdot dz_1}_{\nabla_{z_1}}$$

$$\nabla_{z_1} = \left( \frac{\partial a_1}{\partial z_1} \right)^T \cdot \nabla_{a_1}$$