# Lecture December 13

## The mathematics of NN

input          hidden layers          output
layer                                 layer

$x_1$ ○        $h_1^{(1)}$ ○          $h_1^L$ ○ →
$x_2$ ○        $h_2^{(1)}$ ○          $h_2^L$ ○ →
$\vdots$       ○          - - -       ○ →
$\vdots$       |                      $\vdots$
$\vdots$       |                      $h_k^L$ ○ →
$x_m$ ○        $h_m^{(1)}$ ○

$$W \in \mathbb{R}^{n \times m}$$

input to $-h-$

$$z(x) = w^T x + b$$

$$\sigma(z) = \underbrace{[q_1(x), q_2(x) \cdots q_m(x)]}_{A}$$

with many hidden layers

$$1 \le \ell \le L$$

an Artificial NN cascades
the operations $\sigma(z)$ multiple
times

$$\sigma_L \left( A_L \left( --- \ \sigma_1 \left( A_1(x) \right) \right) \right)$$

consider a simple NN in which $w$ and $b$ are scalars $L = 2$

$$f(x; \Theta) = \sigma_2 \left( w_2 \, \sigma_1 \left( w_1 x + b_1 \right) + b_2 \right)$$

$$a_1 = \sigma_1 \left( w_1 x + b_1 \right)$$

input to $L$ is $a_1 w_2 + b_2$

partial derivatives wrt $w_1$

$$\partial_{w_1} f(x; \Theta) = \sigma_2' \left( w_2 \, \sigma_1 \left( w_1 x + b_1 \right) + b_2 \right) \times w_2 \, \sigma_1' \left( w_1 x + b_1 \right) x$$

with $L$- layers

$$\partial_{w_1} f(x; \Theta) = \left[ \prod_{l=2}^{L} w_l \right] \times \left[ \prod_{l=1}^{2} \sigma_l' (z_l) \right] x$$

$$Z_\ell = A_\ell \left( \sigma_{\ell-1} \left( A_{\ell-1} \left( \cdots \sigma_1 (A_1 x) \right) \right) \right)$$

if $\sigma_\ell$ is the sigmoid
function (or $\tanh$)



$\sigma'_\ell (x)$ will be small if
$|x| >> 0 \implies$ vanishing
gradients

## Typical activation functions

- ReLU function
$$\sigma(z) = \max(0, z)$$
$$\sigma'(z) = 1 \quad \text{for } z > 0$$

- Leaky ReLU
$$\quad \sim ) \quad \begin{cases} \alpha \cdot z & z < 0 \end{cases}$$

$$\sigma(z) = \begin{cases} \\ z & \geq 0 \end{cases}$$

$\sigma(z)$

$\alpha = 0.01$

$z$

$$\sigma'(z) = \begin{cases} \alpha & z < 0 \\ 1 & z \geq 0 \end{cases}$$

$$ELU = \sigma(z) = \begin{cases} \alpha(e^z - 1) & z \leq 0 \\ z & \text{for } > 0 \end{cases}$$

$$\sigma'(z) = \begin{cases} \alpha e^z & z \leq 0 \\ 1 & \text{for } z > 0 \end{cases}$$

$$\sigma(z) = \tanh z \qquad z \in (-1, 1)$$
$$\sigma'(z) = 1 - (\tanh z)^2$$

$$\sigma(z) = \frac{1}{(1 + e^{-z}}$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z)))$$