# FYSS429, FEB 1, 2023

**Basics of NWs**

## Architecture (Model)
- \# hidden layers
- \# nodes/neurons/units
- \# activation function

### Cost function, optimization + regularization

- Cost function

Regression $MSE = \frac{1}{m} \|(t-y)\|_e^2$

$\|x\|_2 = \sqrt{\sum_i x_i^2}$

target

model output

$y = y(\Theta; x)$   our Model

$x =$ input data

$t = f(x)$

$\Theta =$ unknown parameter of our model

cost function $C(\Theta)$

$\hat{\Theta} = \underset{\Theta \in \mathbb{R}^m}{\arg\min} \; C(\Theta)$

- Optimization
  - Gradient descent (GD)
  - GD with momentum

- stochastic GD (SGD)

- algorithms for learning rate
  - Adagrad
  - RMSprop
  - ADAM
  - - - - -

- Regularization

$$\ell_1 \sim \lambda \|\Theta\|_1$$
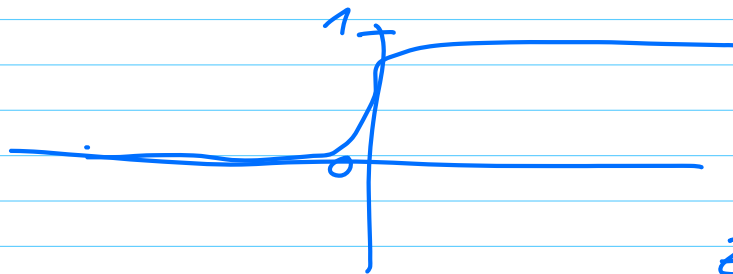
$$\ell_2 \sim \lambda \|\Theta\|_2$$

Cybenko, 1989

Let $\sigma$ be any continuous sigmoidal function

$$\sigma(z) \Rightarrow \begin{cases} 1 & \text{as } z \Rightarrow \infty \\ 0 & \text{as } z \Rightarrow -\infty \end{cases}$$

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



Given a function $F \in [0,1]^d$ and $\varepsilon > 0$, there is a one-hidden layer network $f(x; \Theta)$
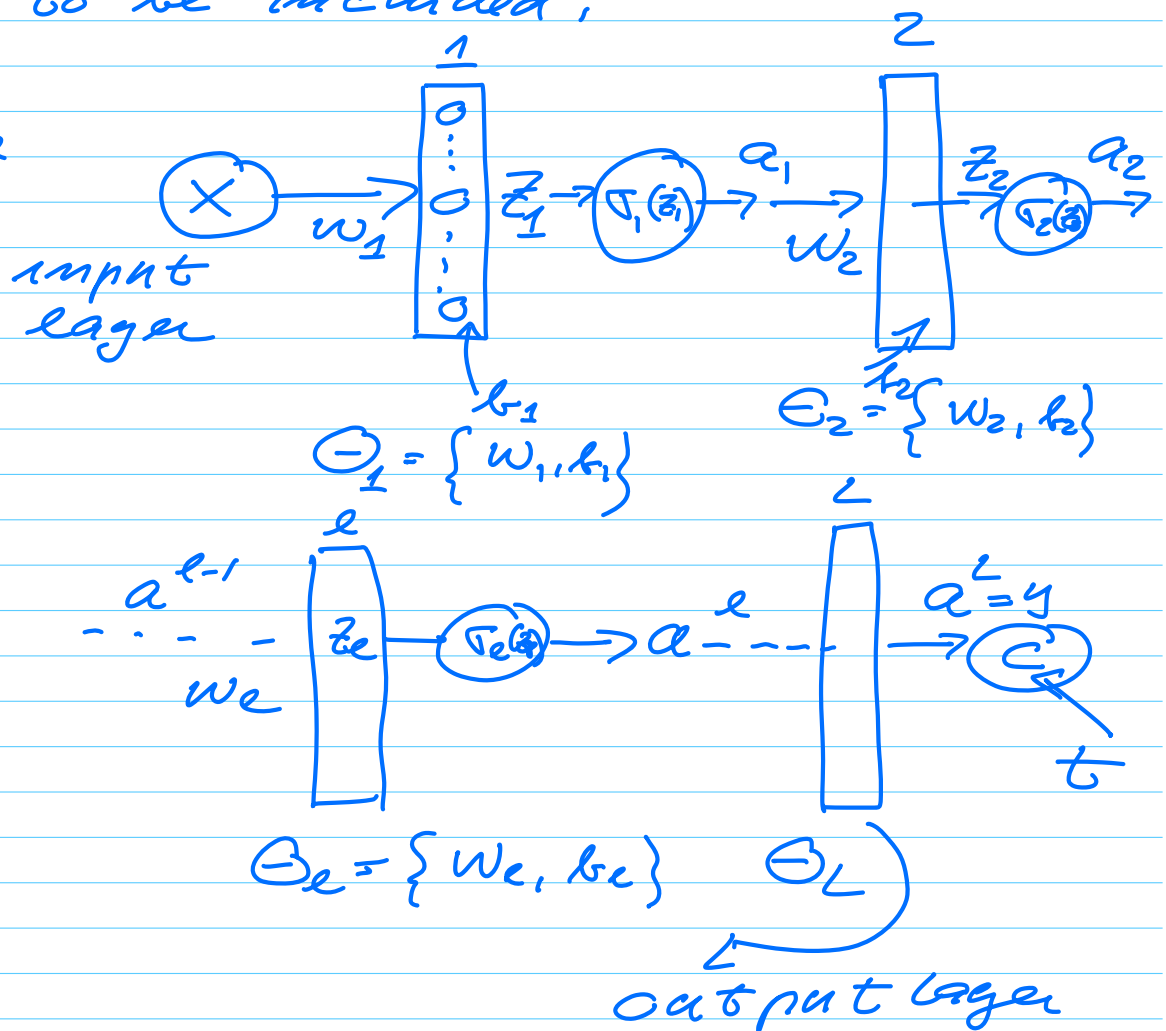
$$\Theta = \{W, b\}$$

$$|f(x; \Theta) - F(x)| < \varepsilon$$

$$\text{for all } x \in C[0,1]^{\alpha}$$

Hornik (1991) extended the theorem to apply to any by letting any non-constant bounded activation function to be included.

<u>NN structure</u>



input layer

$z_1$   $\sigma_1(z_1)$   $a_1$   $W_2$   $z_2$   $\sigma_2(z_2)$   $a_2$

$b_1$

$\Theta_1 = \{W_1, b_1\}$

$\Theta_2 = \{W_2, b_2\}$

$a^{\ell-1}$   $z_\ell$   $\sigma_\ell(z_\ell)$   $a^\ell$   $a^L = y$

$W_\ell$

$t$

$\Theta_\ell = \{W_\ell, b_\ell\}$   $\Theta_L$

output layer

$$\hat{\Theta} = \arg\min_{\Theta} C(\Theta)$$

$$C(\Theta) = \frac{1}{2}\left(t - a^L(x;\Theta_L)\right)^2$$

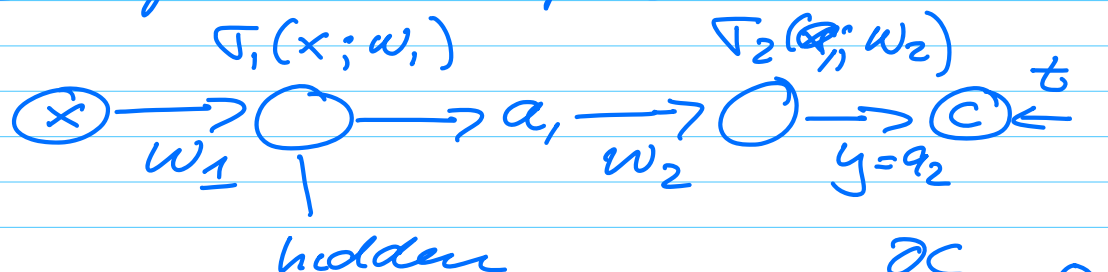$$= \frac{1}{2}\left(t - y(x;\Theta_L)\right)^2 \quad (MSE)$$

Back propagation algo:

$$\frac{\partial C(\Theta)}{\partial \Theta_L} = 0 =$$

$$-\left(t - y(x;\Theta_L)\right)\frac{\partial y}{\partial \Theta_L}$$

$$= -\left(t - a_L(x;\Theta_L)\right)\frac{\partial a_L}{\partial \Theta_L}$$

Simple example



$$a_1 = \nabla_1(x;w_1)$$

$$\frac{\partial C}{\partial w_2} = 0$$

hidden

$$a_1 = \nabla_1(x;w_1)$$

$$y = a_2 = \nabla_2(a_1;w_2) =$$

$$\nabla_2 \left( \nabla_1 (x; w_1), w_2 \right) = a_2$$

$$C = \frac{1}{2}(t - y)^2 = \frac{1}{2}(t - a_2)^2$$

$$\frac{\partial C}{\partial w_1} = ?$$

$$\frac{\partial C}{\partial w_2} = ?$$

activation functions:

$$\nabla_1 (x; w_1) = x \cdot w_1 = a_1$$

$$\nabla_2 (a_1; w_2) = a_1 w_2 =$$
$$w_2 \nabla_1 (x; w_1)$$
$$= y = a_2$$

$$\frac{\partial C}{\partial w_1} = \boxed{-(t-y)} \frac{dy}{\partial a_1} \frac{\partial a_1}{\partial w_1} = \frac{\partial C}{\partial a_2} \frac{\partial a_2}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

$$\frac{\partial C}{\partial w_2} = -(t-y) \frac{\partial y}{\partial w_2} = -(t-y) \frac{\partial a_2}{\partial w_2}$$

$$\frac{\partial C}{\partial w_1} = -(t-y) \times w_2$$

$$\frac{\partial C}{\partial w_2} = -(t-y)w_1 x$$

$$w_1^{(k+1)} \leftarrow w_1^{(k)} - \eta \frac{\partial C}{\partial w_1}\bigg|_{w_1 = w_1^{(k)}}$$

$$w_2^{(k+1)} \leftarrow w_2^{(k)} - \eta \frac{\partial C}{\partial w_2}\bigg|_{w_2 = w_2^{(k)}}$$

$$\frac{\partial C}{\partial \theta_L} = \frac{\partial}{\partial \theta_L}\left( C(t, a_L(\theta_L; x)) \right)$$

$$\text{if } C = MSE = \frac{1}{2}\|t - a_L\|_2^2$$

$$\frac{\partial C}{\partial \theta_L} = -(t - a_L) \frac{\partial a_L}{\partial \theta_L}$$

$$\frac{\partial C}{\partial \theta_{L-1}} = \frac{\partial C}{\partial a_L} \frac{\partial a_L}{\partial \theta_{L-1}}$$

$$\frac{\partial C}{\partial \theta_{L-2}} = \frac{\partial C}{\partial a_L} \frac{\partial a^L}{\partial a_{L-1}} \frac{\partial a_{L-1}}{\partial \theta_{L-2}}$$

$$\frac{\partial C}{\partial \theta_\ell} = \boxed{\frac{\partial C}{\partial a_L}} \frac{\partial a_L}{\partial a_{L-1}} \frac{\partial a_{L-1}}{\partial a_{L-2}} \cdots$$

$$\cdots \quad \frac{\partial a_{\ell+2}}{\partial a_{\ell+1}} \quad \frac{\partial a_{\ell+1}}{\partial B_\ell}$$

Automatic Diff | in python import autograd
JAX

## Example

$$f(x) = \sqrt{x^2 + \exp(x^2)}$$

$$\frac{df}{dx} = \frac{x\,(1 + \exp(x^2))}{\sqrt{x^2 + \exp(x^2)}}$$

$f(x)$ in a brute force way

$$x^2 = x \cdot x = 1 \text{ FLOP}$$
$$\exp(x^2) = \exp(x \cdot x) = 2 \text{ FLOPS}$$
$$x^2 + \exp(x') = 1 \text{ FLOP}$$
$$\text{SQRT} = -1$$
$$\overline{\phantom{xxxxxxxxxx}}$$
$$5 \text{ FLOP.}$$

$$a = x^2 \qquad \exp(a) \Rightarrow$$
$$4 \text{ FLOP.}$$

## Derivative

Numerator : 4 FLops
Denominator : 5 FLopT
+ Division ; 1 Flop
$$\overline{\phantom{xxxxxxxxxx}}$$
10 FLops

$$\frac{df}{dx} \simeq \frac{f(x + \Delta x) - f(x)}{\Delta x}$$

# Automatic diff:
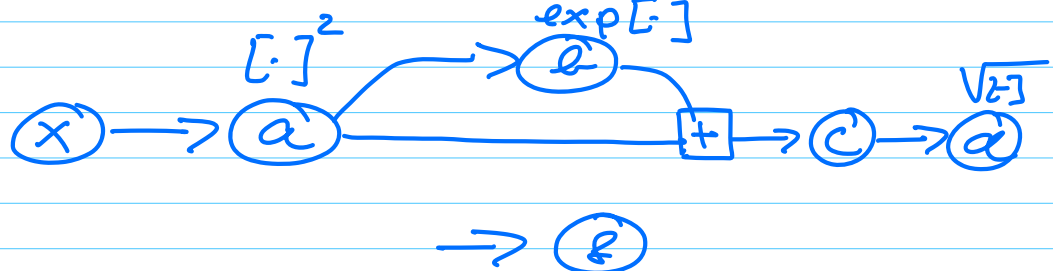
$$f(x) = \sqrt{x^2 + \exp(x^2)}$$

- Forward mode
- Reverse mode (AKA Back prop)

$$a = x^2 \qquad b = \exp(x^2) = \exp(a)$$

$$c = a + b$$

$$d = \sqrt{c} = f(x)$$



$$\frac{df}{dx} = ?$$

$$\frac{da}{dx} = 2x \qquad \frac{db}{dx} = \frac{db}{da}\frac{da}{dx}$$

$$= 2x \exp(x^2)$$