# Artificial Neural Networks

Machine Learning for Data Science 1

CENTRAL IDEA :

- extract derived features as linear combinations of inputs
- model the internal target features through non-linear transformations
- recursive and ~~repeat~~  1, >2 ≡ deep model

1943 : models of the true NN

1940s : plasticity $\equiv$ learning

1953 : perceptron , failed idea , looked very promising

1965 : Idea of many layers

1973 : backpropropogation ||

1975 : $\rightarrow$

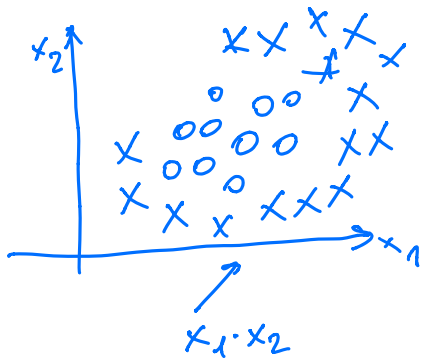1989: LeCun : hand-written digits $\}$ deep networks convolution

1992 : max pooling , 3D

2006 : Hinton : Boltzman NN

2009 - 2012 : ANN major competitions

2012 : deep learning , images, text, ---- more data comp. power

MOTIVATION

- can we learn "hard concepts"
- Interactions

1000 features

$$\rightarrow \frac{1000 \cdot 999}{2} \text{ two-interaction}$$

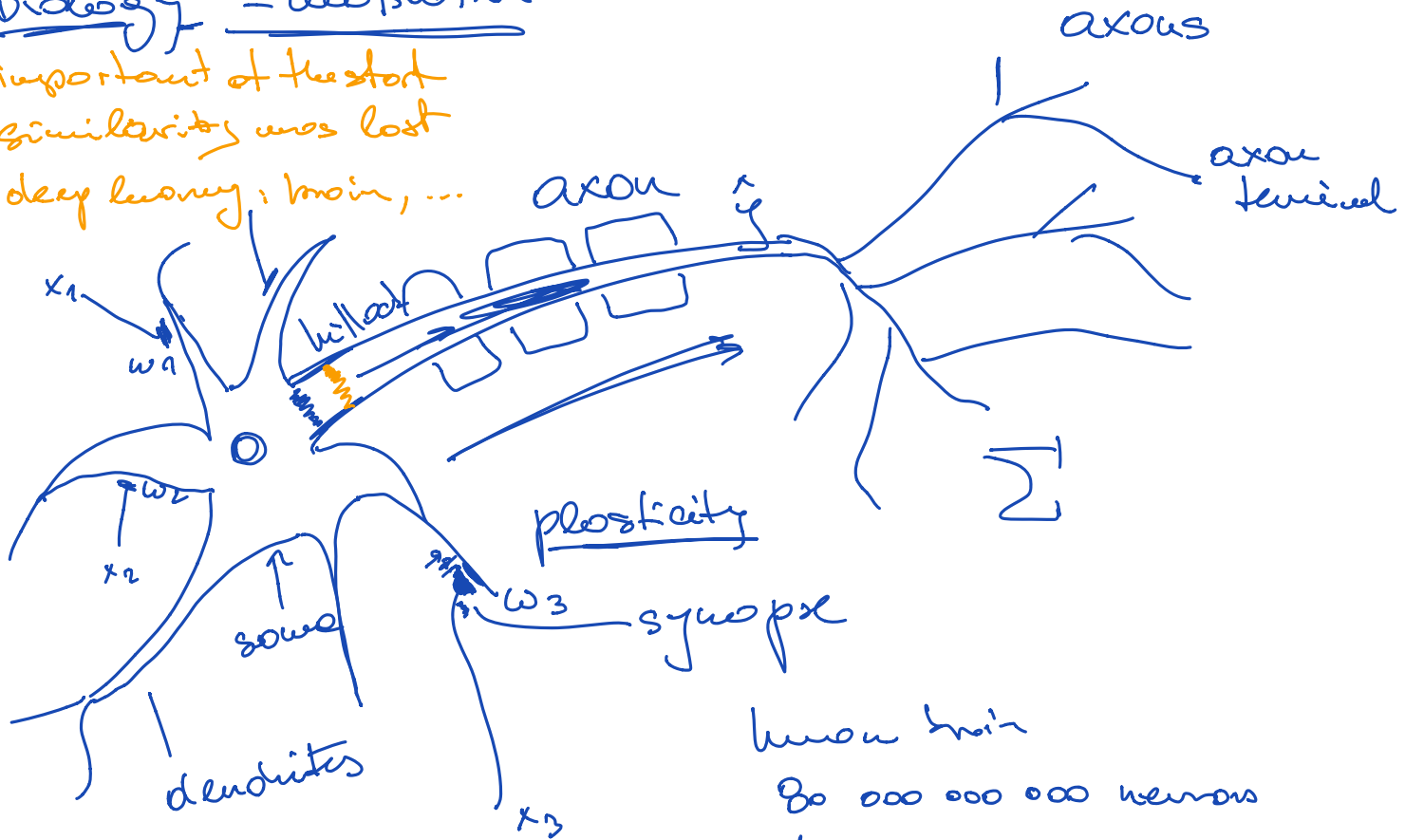$$\frac{1000 \cdot 999 \cdot 998}{3 \cdot 2} \text{ three interactions}$$

complexity
needs a lot of data
overfitting
explanation

a modelling framework
that incorporates all
these interactions

potentially

$x_1 \cdot x_2$

# Biology — motivation

important at the start
similarity was lost

deep learning, brain, ...

axons

axon ŷ

hillock

$x_1$

$w_1$

$w_2$

$x_2$

soma

$w_3$

plasticity

synapse

dendrites

$x_3$

axon terminal

$\sum$

synapses are slow
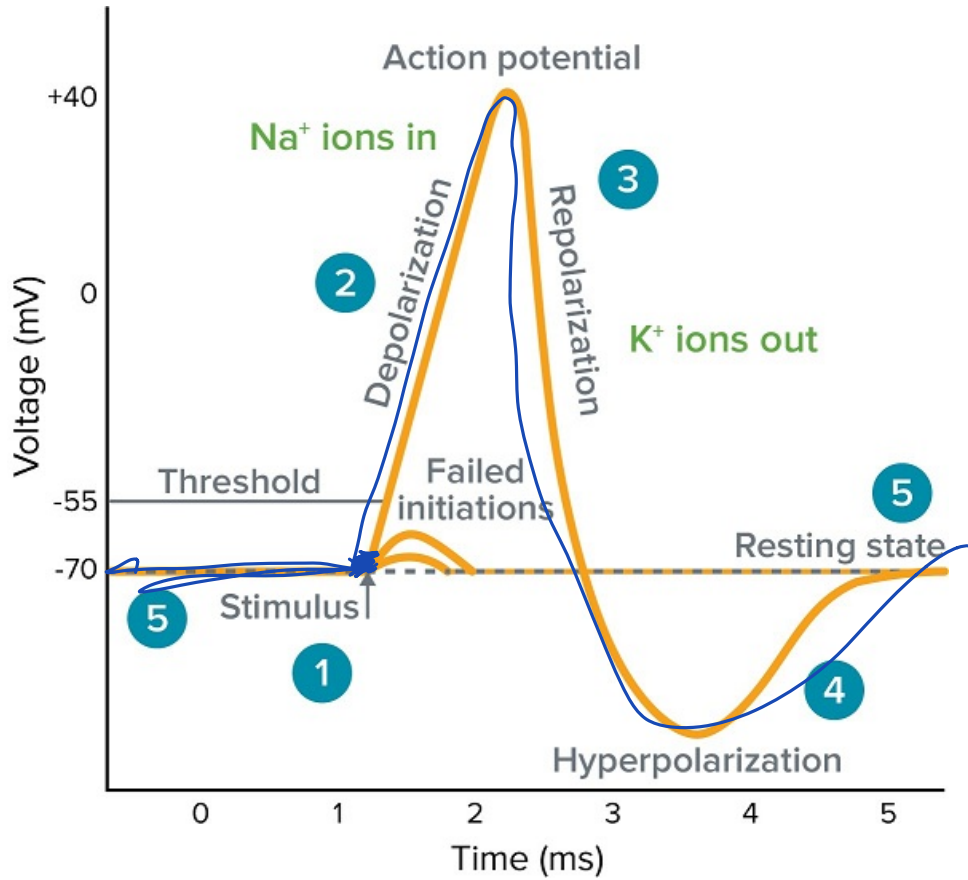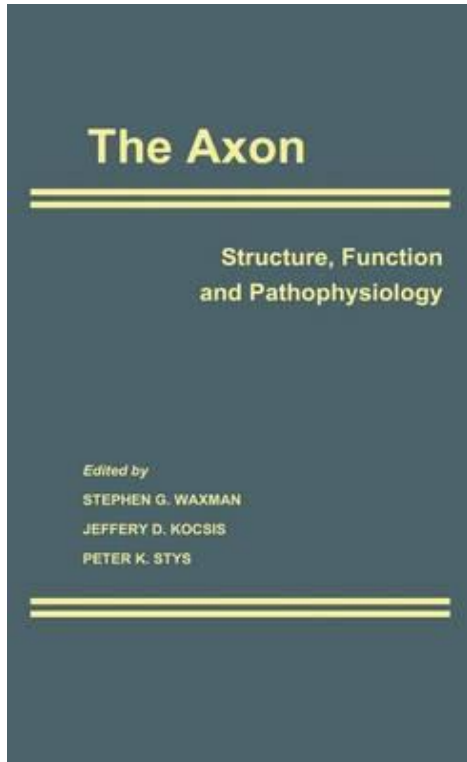
everything is ||

human brain
80 000 000 000 neurons
10 000 synapses
---
$10^{15}$ connections

# Modularity of the Brain

- different regions = different tasks
- anatomically similar
- damage in one region
  $\Rightarrow$ a different region can take over

# The Axon

## Structure, Function and Pathophysiology

Edited by

STEPHEN G. WAXMAN

JEFFERY D. KOCSIS

PETER K. STYS

---

# 1 | Electrical activity of nerve: The background up to 1952

SIR ANDREW HUXLEY

My interest in physiology, and in the physiology of nerve in particular, dates from the autumn of 1935, when I went up to Trinity College, Cambridge, as an undergraduate. I was expecting to specialize in physics, in which I had been very well taught at school, but the rules of the Natural Sciences Tripos ("tripos" is a Cambridge word for courses leading to a first degree) required me to take a third experimental science as well as the physics, chemistry, and mathematics that were the obvious choices. I picked physiology on the advice of a friend a few years older than myself who told me that it was a lively subject in which even the initial courses included material recently discovered or even still controversial, unlike the courses in physics, which included nothing that had not been cut and dried for decades. I was inspired to switch to physiology as my final-year specialty subject largely by my teachers W. A. H. (William) Rushton and F. J. W. (Jack) Roughton and by personal contacts with Glenn Millikan (son of R. A. Millikan of the oil-drop experiment; too little known on account of his death in 1946 in a climbing accident) and Alan Hodgkin, all Fellows of Trinity College working in the physiology laboratory. E. D. Adrian (later Lord Adrian, Master of Trinity College and President of the Royal Society) was also a Fellow of Trinity College, but I hardly came across him until my final undergraduate year because he was a research professor of the Royal Society, taking little part in undergraduate teaching, until 1937; in that year he became head of the Department and in my final year he lectured to us on the central nervous system.

I hope that my account of the ideas then current about nerve conduction and of developments up to 1952 is not too heavily biased by my Cambridge background.

## EXCITATION OF NERVE

Our first-year lectures on nerve were given by William Rushton. We were, of course, taught the elementary facts about excitation of nerve, mostly established in the mid-19th century in Germany by experiments on the sciatic nerve of the frog with the gastrocnemius muscle attached to indicate by its contraction whether the motor nerve fibers had been activated: the impulse arises at the point where a stimulating current of short duration leaves the nerve (the cathode) and it travels in both directions. If a direct current of fairly long duration is used, an impulse may be set up both at the cathode at the start of the current and at the anode when the current is terminated (anode break excitation). There is no response if the strength of the stimulus is below a well-defined critical level (the "threshold"), and the response increases with stimulus strength up to a maximum; the impulse is accompanied by a wave of electrical negativity passing along the surface of the nerve. A second stimulus is ineffective if it follows a maximal stimulus within a certain time interval (the "absolute refractory period," roughly equal to the duration of the propagated electric change), and this is followed by a "relative refractory period" in which the threshold is higher than when the nerve is fully rested. The threshold value of current strength varies inversely with its duration, the product of these quantities approaching a finite limit as the duration is reduced toward zero.

## THE ALL-OR-NONE "LAW"

At the turn of the century, it had been debated whether the gradation of response with strength of stimulus was solely a matter of the number of fibers within the nerve trunk being activated, or whether the impulse in an individual fiber could vary with the strength of the stimulus. The former alternative was found to be correct: the invariant "all-or-none" character of the propagated response of individual motor nerve fibers, and of individual fibers of skeletal muscle, was well established in the first decade of this century by Keith Lucas (1905, 1909) (another Fellow of Trinity College, and, like Millikan, too little known on account of his early death in a flying accident during World War I), using the twitch of a muscle fiber or a motor unit as the indication of activity. It was recognized that the energy dissipated by
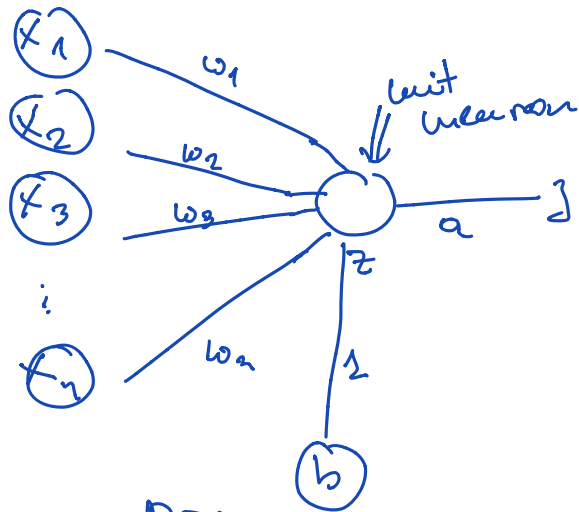
3

---

**Format**: Hardback | 708 pages

**Dimensions**: 225 x 289 x 39mm | 2,278g

**Publication date**: 30 Mar 1995

**Publisher**: Oxford University Press Inc

**Publication City/Country**: New York, United States

**Language**: English

**Edition Statement**: New

**Illustrations note**: halftones, line figures and tables

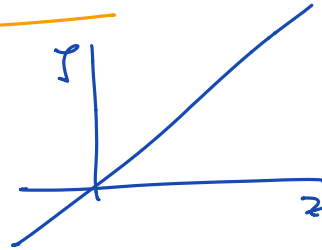**ISBN10**: 0195082931

**ISBN13**: 9780195082937

$x_1$ $w_1$

$x_2$ $w_2$

$x_3$ $w_3$

$\vdots$

$x_n$ $w_n$

unit neuron

$z$ $a$ $y$

$1$ $b$

**Linear neuron**

$$z = \sum w_i x_i + b$$

$$y = \sigma(z) \qquad \sigma(z) = z$$

Linear Unit

**Rectified Linear**

RELU

**Sigmoid neurons**

$$y = \frac{1}{1 + e^{-z}}$$

# Perceptrons

1960s : Frank Rosenblatt

how to learn weigh vectors

Weigh-space

$$z = \omega^T x$$



good

bad

Decision lines

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$$y = x_1 \equiv x_2$$

initialize $\omega$

repeat

randomly choose $x$ from training set

if $h(z) \neq y$:

if $h(z) = 0$:

$$\omega \leftarrow \omega + x$$

else:

$$\omega \leftarrow \omega - x$$

# Some conventions

- units of ANN, activation $\in 0..1$
- output of ANN, <u>real values</u> $\Rightarrow$ regression
  <u>probabilities</u>

$$
\begin{array}{ll}
0 & cat \\
0 & mouse \\
0 & dog \\
0 & \vdots \\
0 & \vdots \\
\end{array}
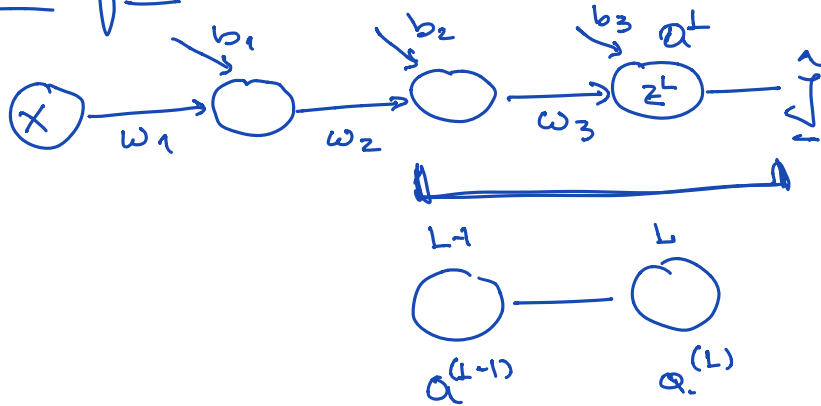\quad \Bigg| \quad \Sigma = 1
$$

$k$

softmax

$$\hat{y}_i(z) = \frac{e^{z_i}}{\sum_{i=1}^{k} e^{z_i}}$$

$k = 2 \equiv$ logistic regression
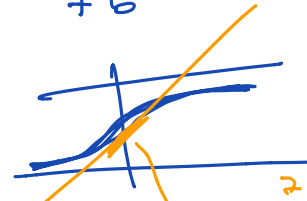
$k > 2$

# example



$$a^{(L)} = \sigma\left(z^{(L)}\right)$$

$$z^{(L)} = \sum_i \omega_i^{(L)} x_i^{(L)} + b^{(L)}$$
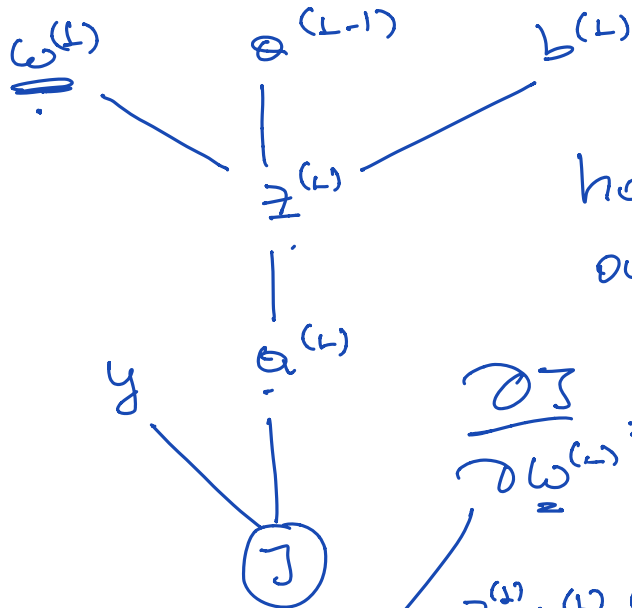
$$J(\omega_1, b_1, \ldots \omega_3, b_3) = \left(a^{(L)} - y\right)^2$$

$$z^{(L)} = \omega^{(L)} a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma\left(z^{(L)}\right)$$

$$\frac{\partial \varphi(z)}{\partial z} = \varphi(z)(1 - \varphi(z))$$

linear part of the sigmoid

$w^{(L)}$

$a^{(L-1)}$

$b^{(L)}$

$z^{(L)}$

$a^{(L)}$
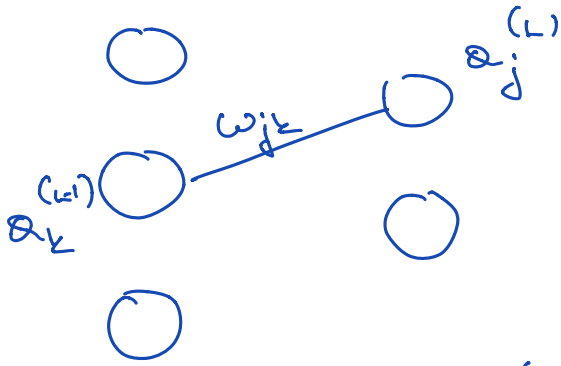
$y$

$J$

how does $J$ depend
on $w^{(L)}$

$$\sum$$
$$J = \left(a^{(L)} - y\right)^2$$

$$\frac{\partial J}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}} \cdot \frac{\partial J}{\partial a^{(L)}}$$

$z = w^{(L)} a^{(L-1)} + b^{(L)} \qquad \sigma(z)(1-\sigma(z)) \qquad 2\left(a^{(L)} - y\right)$

$a^{(L)}$ \qquad\qquad\qquad\qquad error

$$\frac{\partial J}{\partial w^{(L)}} = a^{(L-1)} \quad \sigma(z)(1-\sigma(z)) \cdot 2\left(a^{(L)} - y\right)$$

$$J = \sum_{j=0}^{n} \left( a_j^{(L)} - y_j \right)^2$$



$$z_j^{(L)} = \sum_i \omega_{ji}^{(L)} a_i^{(L-1)} + b^{(L)}$$

$$a_j^{(L)} = \sigma\left( z_j^{(L)} \right)$$

$$J = \sum_i \left( a_j^{(L)} - y_j \right)^2$$

$$\frac{\partial J}{\partial \omega_{ji}^{(L)}} = \sum \frac{\partial z_j^{(L)}}{\partial \omega_{jk}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial J}{\partial a_j^{(L)}}$$

$$\frac{\partial J}{\partial a_k^{(L-1)}} = \sum \frac{\partial z^{(L)}}{\partial a_k^{(L-1)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial J}{\partial a_j^{(L)}}$$

$$L-1 \leftarrow$$

$$\frac{\partial J}{\partial a_k^{(L-2)}} \cdots \int \left( \frac{\partial J}{\partial a_k^{(L-1)}} \right) \cdots$$

$$\underline{\underline{X}} = \left[\equiv\right]^{u_1} \underset{m}{} \quad , \quad \overrightarrow{X}' = \left[1\,|\,\equiv\right]^{u_1+1} \underset{m\ 1}{}$$

$$\underline{X}^{(1)} = \underline{X}'$$

$$Z^{(2)} = \underline{X}^{(1)} \, W^{(2)}$$

$$\underset{m\times n_2}{} \quad \underset{m\times n_1}{} \quad \underset{n_1\times n_2}{}$$

$$+ \frac{\otimes}{m} \sum_{\ell} \sum_{ij} W_{ij}^{(k)}$$

$$\underline{A}^{(2)} = \sigma(Z^{(2)})$$

$$\underline{X}^{(\ell)} = \sigma\left(A^{(\ell-1)} \, \omega^{(\ell)}\right)$$

$$\frac{\partial J}{\partial \omega^{(L)}} = \frac{\partial Z^{(L)}}{\partial \omega^{(L)}} \cdot \frac{\partial a^{(L)}}{\partial Z^{(L)}} \cdot \frac{\partial J}{\partial a^{(L)}} \cdots \frac{\partial J}{\partial \varrho^{(L-1)}} \cdot \left(\frac{\partial \varrho}{\partial a^{(k)}}\right)$$

$$\underset{m\times n_L}{d^{(L)}} = A^{(L)}\left(1 - A^{(L)}\right) \circ \left(A^{(L)} - \overrightarrow{Y}\right) -$$

$$\underline{J}^{(L)} = \frac{1}{m}\left(A^{(L-1)}\right)^T \times d^{(L)} \quad -$$

# Learning

  Optimization : (stochastic) gradient descent

  generalization : tricks $\longrightarrow$ mini-batch , 250
  adaptive learning rate
  momentum

  - regularization
    "weight -decay"
      0.99

  - weight sharing
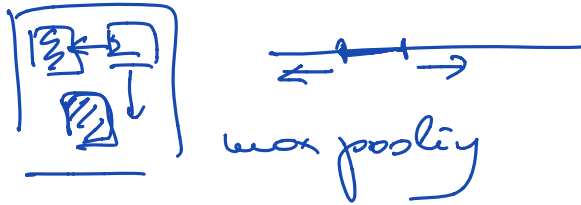  - early stopping | $\rightarrow$ validation set
  - ensambling
  - drop-out
  = pre-training || transfer learning

# deep learning

- hidden layers $\geq 2$

  - convolutional NN

    

    max pooling

  - recurrent NN

  - long short-term memory units

  - transfer learning

  - encoders