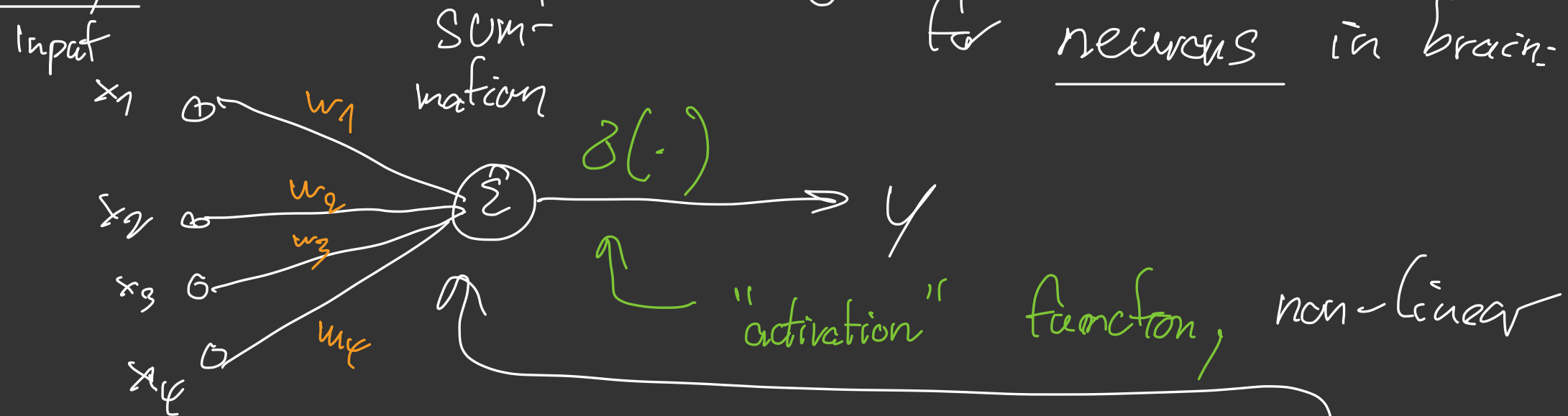


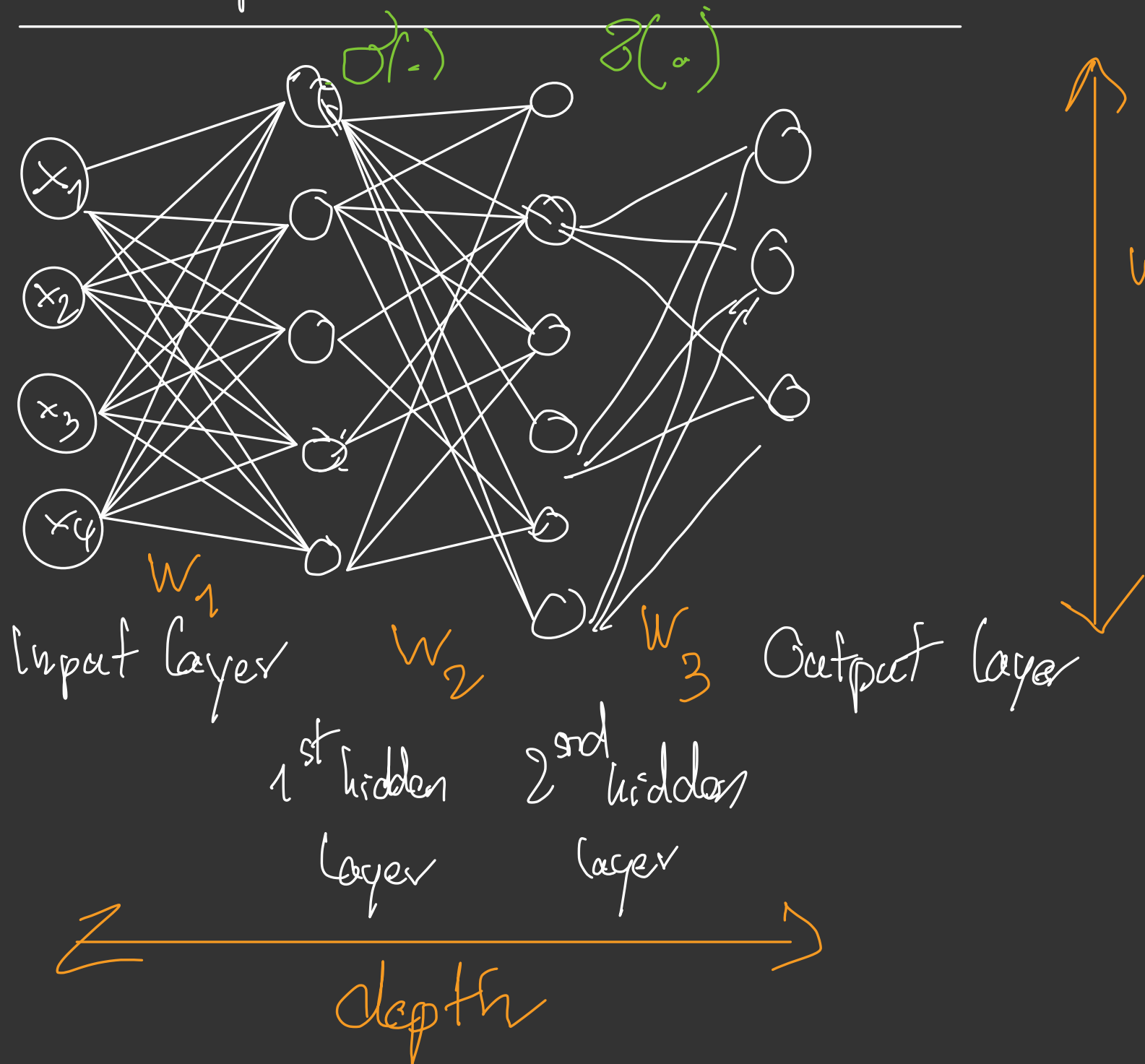
# Artificial Neural Networks

(Selective) Little History: [McCulloch, Pitts '43]: Mathematical model for neurons in brain.



- 1950's: Perception; 1<sup>st</sup> numerical scheme to use collection of ) to do classification
- 1980's: "Deeper" networks, successful training via **backpropagation** (training algorithm)
- Since ~2007:
  - ▷ Efficient training ( $\hat{=}$  determination of weights ( $w_i$ ) from training data set) becomes possible for deeper networks  $\rightarrow$  regularization
  - ▷ Advances in computing technology: Use GPUs
  - ▷ Outperformance of traditional learning methods (such <sup>as</sup> support vector machines)

# A template for neural networks



A typical hypothesis space  $\mathcal{F}$  associated to an ANN:

$$\mathcal{F} = \{h: X \rightarrow Y\}$$

$$h(x) = \sigma(W_L(\sigma(W_{L-1}(\sigma(\dots \sigma(W_1(x))\dots))),$$

$\sigma$  non-linear, acts component wise

$$W_l: \mathbb{R}^{N_{l-1}} \rightarrow \mathbb{R}^{N_l} \quad \left. \begin{array}{l} \text{(affine) Linear} \\ \forall l=1, \dots, L \end{array} \right\}$$

$N_0$ : dimension of input layer

$N_i$ :  $i$ -th hidden layer

$N_L$ : of output layer

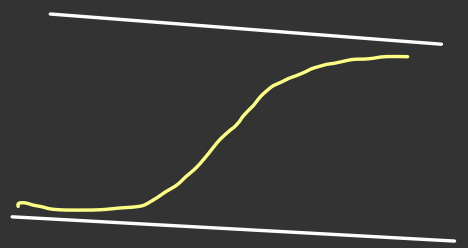
Example: • (Multiclass) Logistic Regression: 1-layer NN ( $L=1$ ),

$$\sigma(\cdot) = \text{softmax}(\cdot)$$

"sigmoidal"

•  $\forall f \downarrow$  with  $L > 1$ :

Multilayer Perceptron



Thm: [Horvick, '89; Cybenko, '91]:

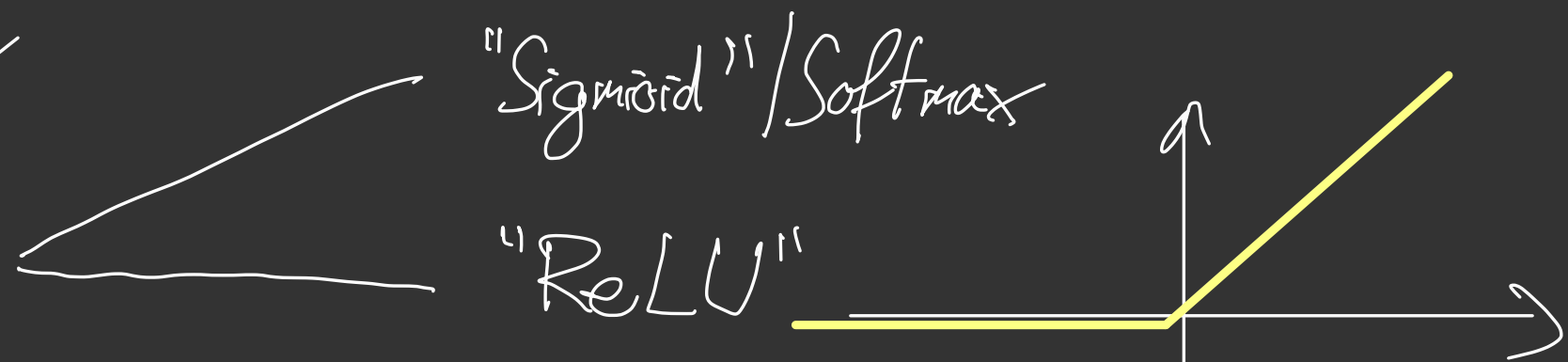
Every measurable function can be approximated by a  $\text{NN}$  with at least  $L=2$  layers (if wide enough).

Practical Q:

▷ How many layers  $L$ , how wide each layer?

▷ Choose which loss fct?

▷ Which activation function?

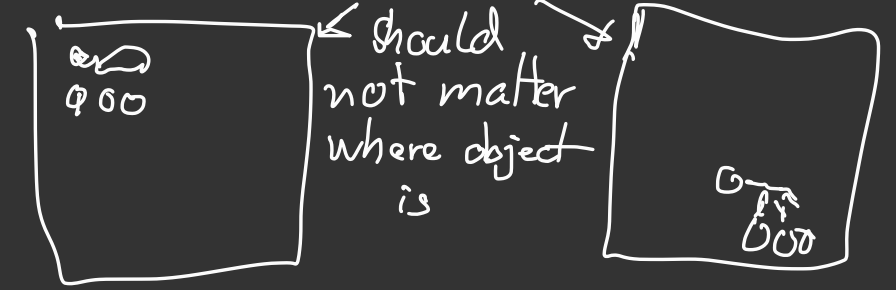


▷ How to "train" the network (i.e., determine weights  $W_e$ )

↳ Some variant of stochastic gradient descent (SGD), implemented

via backpropagation (i.e., a very smart implementation of the chain rule from calculus)

## Modifications / Variants:



▷ Constrain weights to have certain properties: **Convolutional Neural Network (CNN)**  
good for imaging problems (enforces spatial invariance)

▷ Fewer connections: → often less overfitting

▷ (Max / Average) Pooling: Reduces nr. of parameters, avoids overfitting

12	20	30	0
8	12	2	0
34	70	37	4
112	100	25	12

→ 2 × 2 Max-Pool →

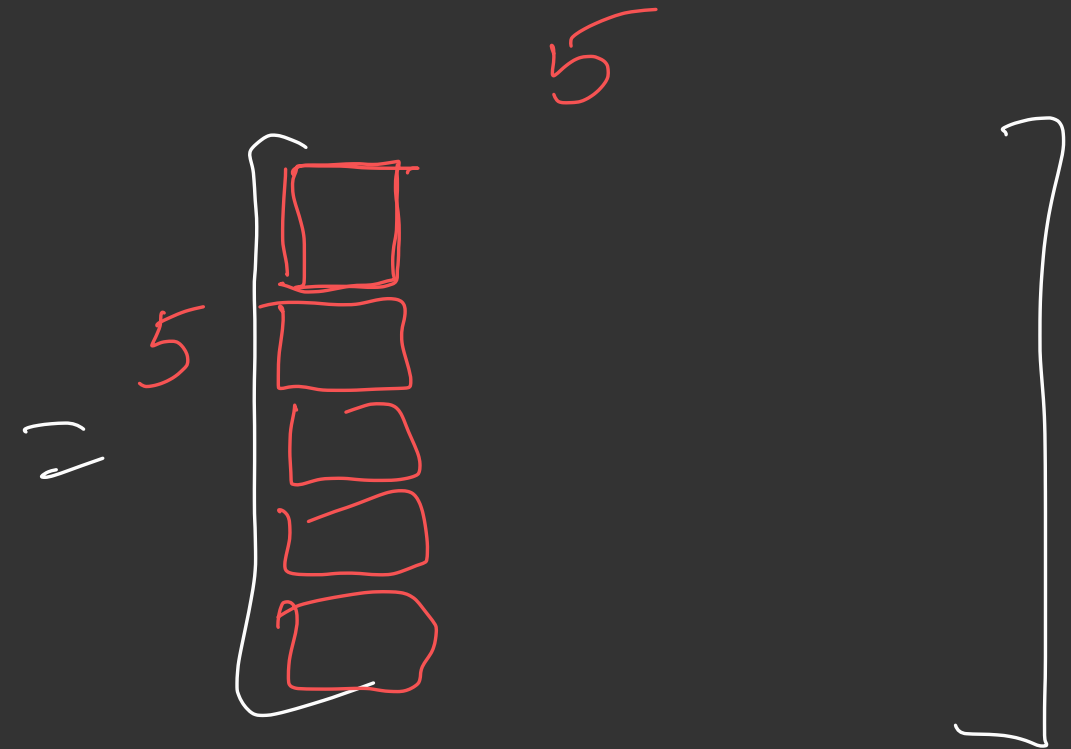
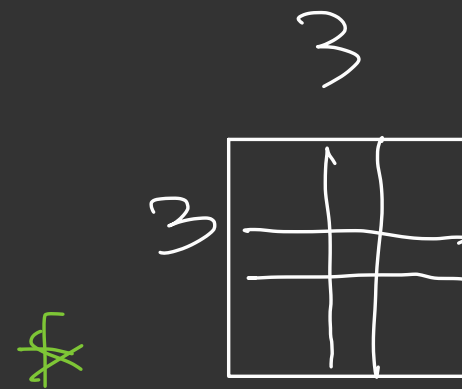
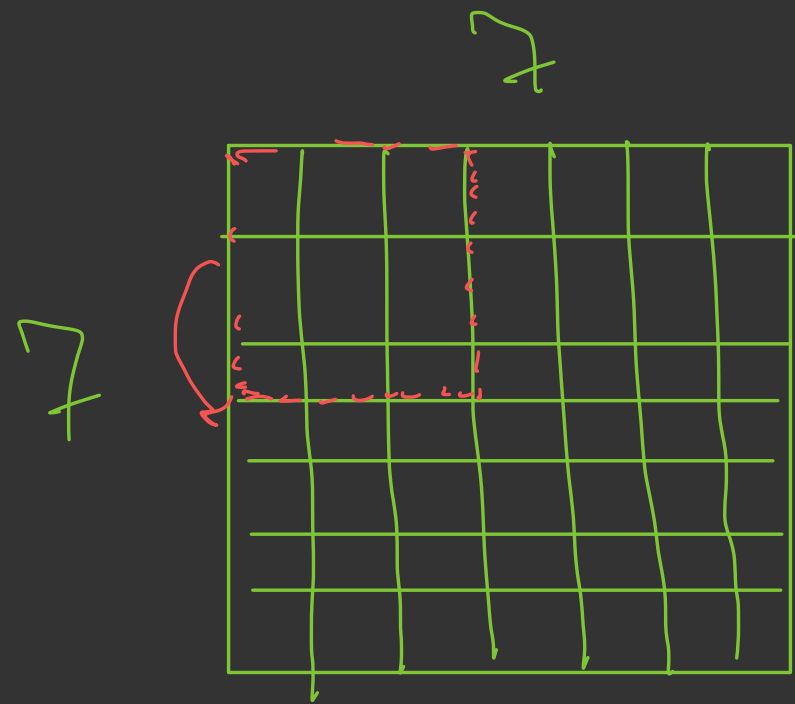
20	30
112	37

▷ Dropout: Drops (at random) connections between layers in training phase

▷ Batch normalization: Normalizes input of a layer,  
↳ faster training, better generalization

Epoch: One pass through entire training data set.

Example: Convolutional 3x3 filter



Reduction from  $7 \times 7 = 49$  parameters

to

$5 \times 5 = 25$  parameters