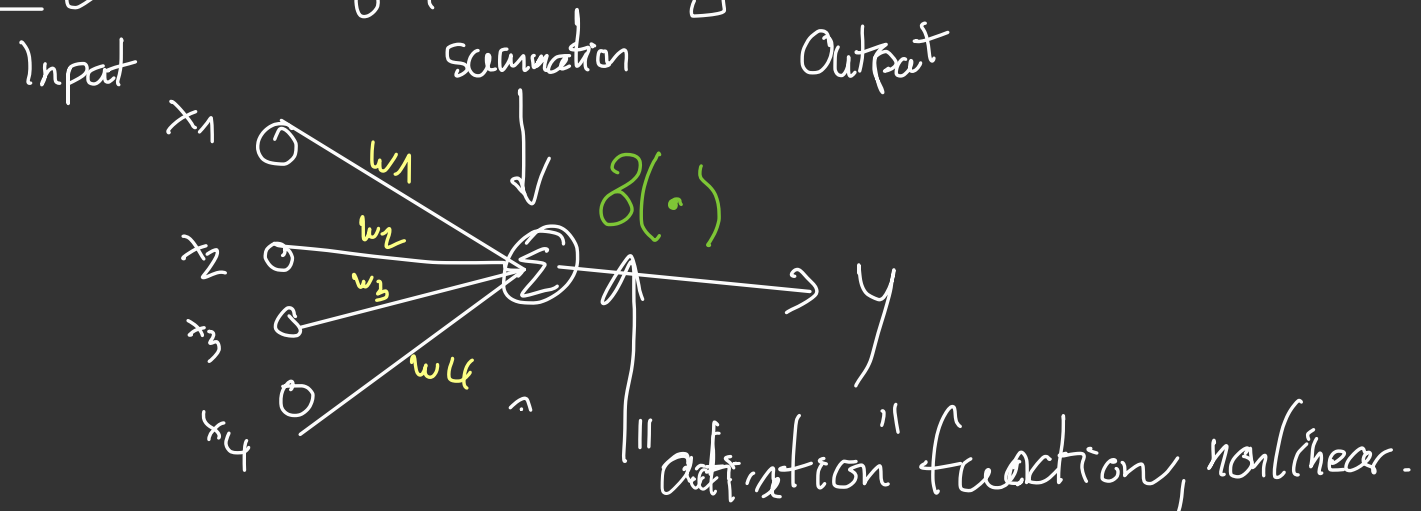


Artificial Neural Networks

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

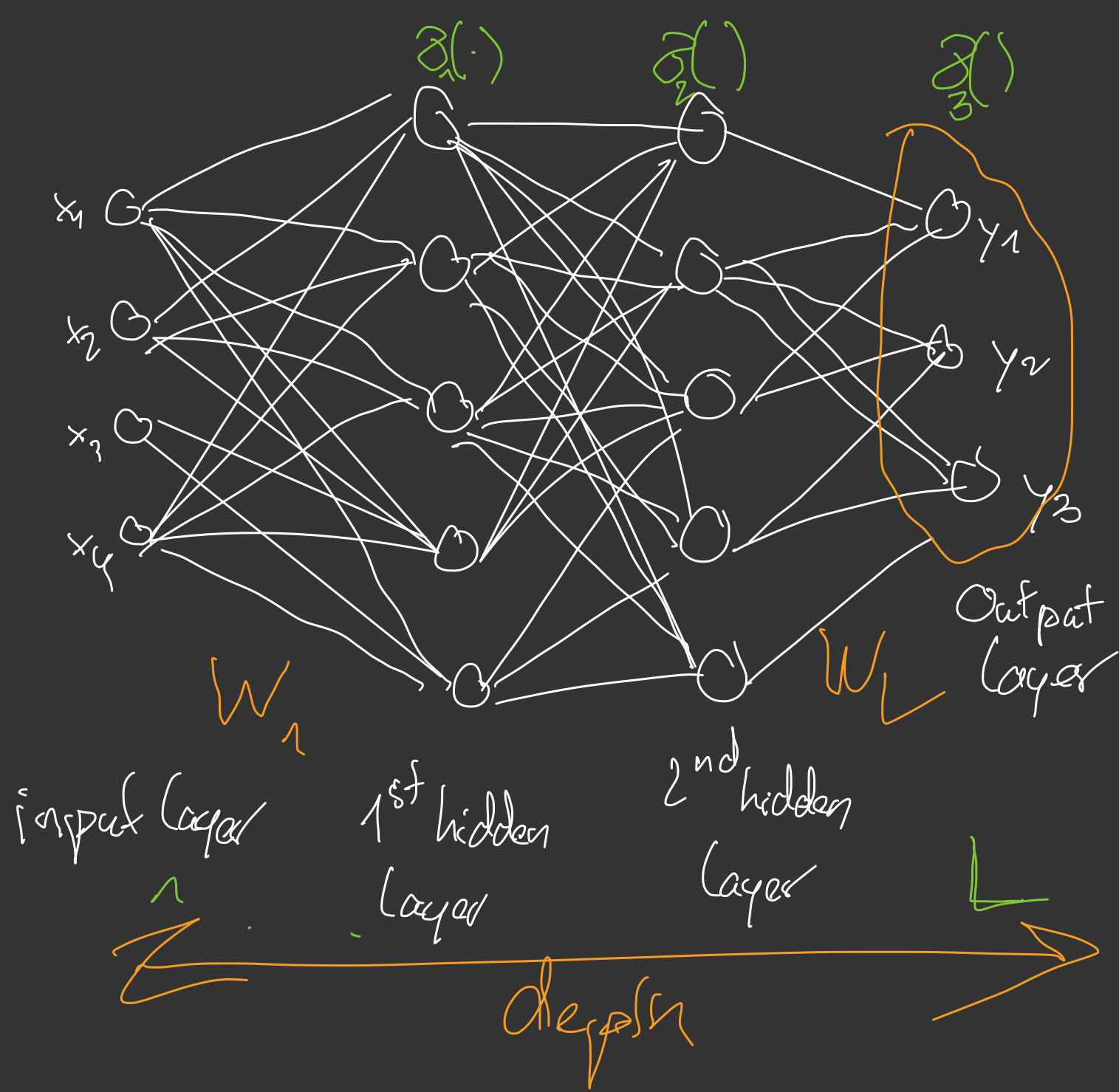


1950's: Perceptron: 1st numerical scheme to use in classification

1980's: "Deep" networks, successful training ... backpropagation

Since ~2007: \triangleright Efficient training, Regularization

- \triangleright Advances in computing technology: Use GPU's
- \triangleright Outperformance of traditional learning methods



A typical hypothesis space \mathcal{F}

with $\mathcal{F}_L = \left\{ h: X \rightarrow Y: \right.$

$$h(x) = \sigma_L \left(W_L \left(\sigma_{L-1} \left(W_{L-1} \left(\dots \left(\sigma_1 \left(W_1(x) \right) \right) \right) \right) \right) \right)$$

$\sigma_1, \dots, \sigma_L$ are non-linear, acts component wise

$W_l: \mathbb{R}^{N_{l-1}} \rightarrow \mathbb{R}^{N_l}$ (affine) linear

If $l = 1, \dots, L$

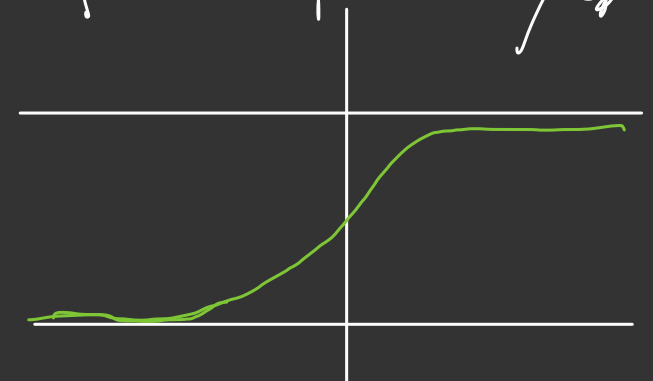
N_0 : dimension of input layer

N_l : — " — l -th hidden layer

Example: Logistic Regression: 1-layer N ($L=1$) N_L : — " — of output layer

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

if $L > 1$: Multilayer Perceptron



Thm: [Hornik, Cybenko]:

Every measurable function can be approximated by a neural network with at least $L=2$ layers (if wide enough).

Practical Question:

▷ How many layers L ?

▷ Which loss function l ?

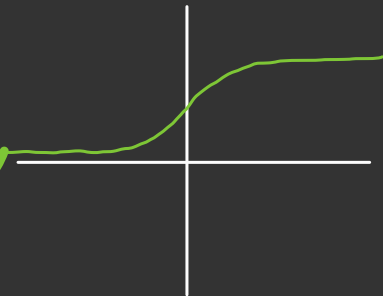
▷ Which activation function?

▷ How to "train"?

• Stochastic Gradient Descent

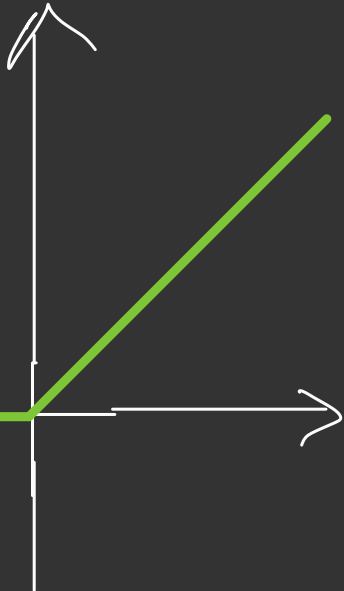
• Adam

Sigmoid



Rectified
Linear Unit

"ReLU"



Modifications / variants:

- Constrain weights to have certain properties: Convolutional Neural Networks (CNN)
good for image problems (enforce spatial invariance).
- Fewer connections \rightarrow less overfitting
- (Max/Average) Pooling:
Reduces spatial sensitivity
- Dropout: Drops (at random) connections between layers in training phase
- Batch normalization: Normalizes input of a layer
 \rightarrow faster learning, better generalization

Epoch: One pass through entire training data,

Learning Rate: Step size of θ (SSD)