

CMC MSU Department of Algorithmic Languages
Samsung Moscow Research Center

Neural Networks for Natural Language Processing

Нейронные сети в задачах автоматической обработки текстов

Lecture 3. Introduction to Neural Networks

Arefyev Nikolay

*CMC MSU Department of Algorithmic Languages &
Samsung Moscow Research Center*

Inspired by the brain

Artificial NNs are inspired by the brain, but are very different (like airplanes are inspired by birds, but are different).

The brain consists of $\sim 10^{11}$ neurons, each connected to $\sim 10^3$ - 10^4 other neurons, a few are connected to receptors.

- Very high level of parallelism
- Neurons communicate using spikes - binary signals (charge)
- Learning algorithm?

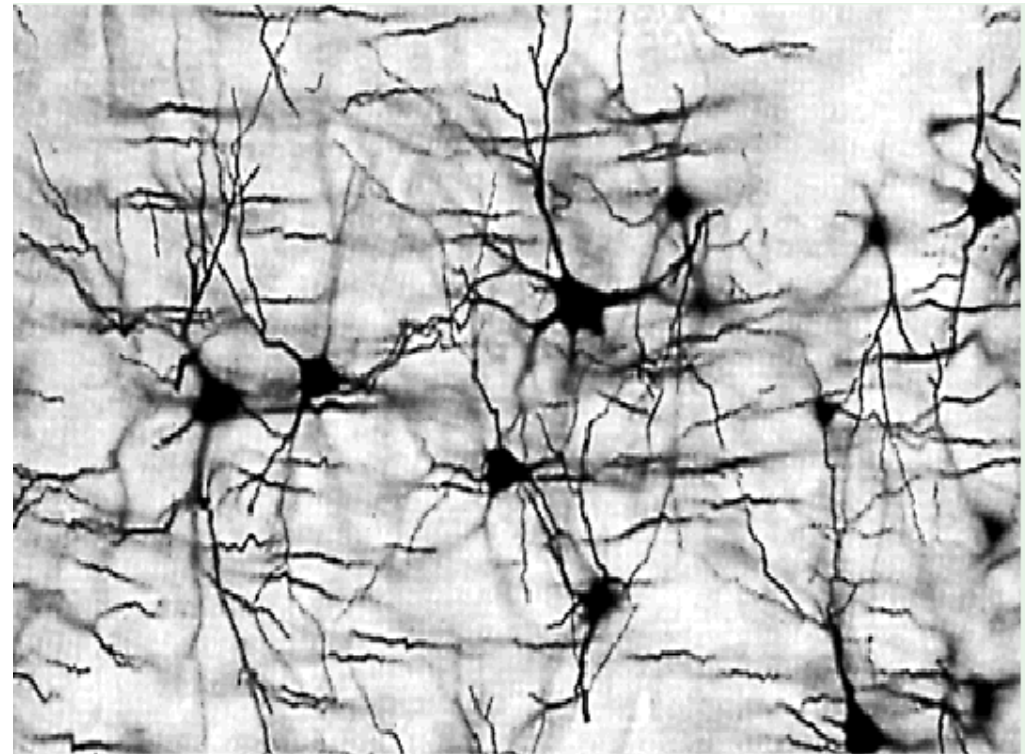
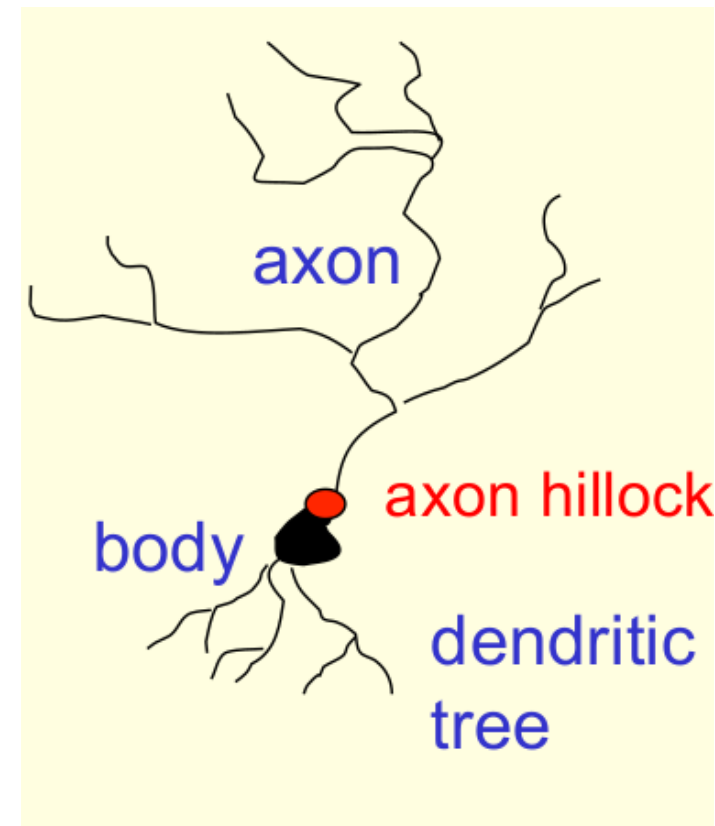


Figure from Yaser Abu-Mostafa. Machine learning video library (<https://work.caltech.edu/library/>) ²

Cortical neuron

Dendrites receive signals from other neurons (inputs), **axon** sends signals to other neurons (outputs). They are connected with special structures – **synapses** ($\sim 10^{14}$).

- Axon hillock generates spike when enough charge is in
- 2 types of synapses: excitatory / inhibitory (increase / decrease probability of a spike)
- Synaptic weights (strength of influence on probability of a spike) are adapted to do computations useful for real-life problems



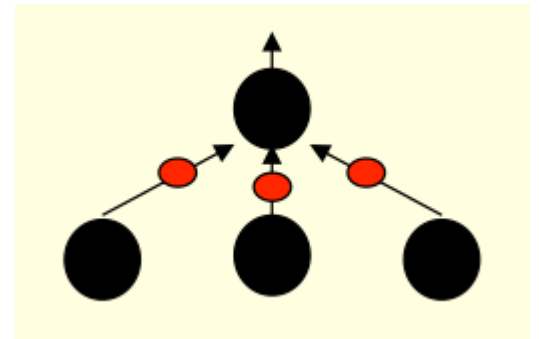
Artificial neurons

Most of them are wrong models of real neurons ...

- For instance, use real values to communicate ... but allow to apply effective algorithms to do useful things.

Artificial neurons first compute weighted sum of inputs (**preactivation**):

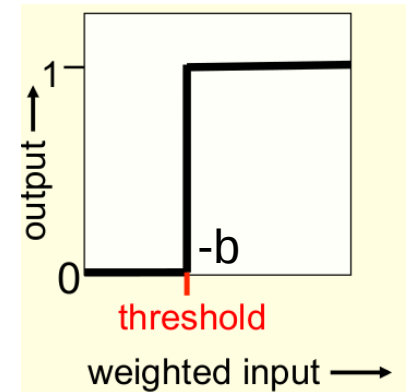
$$z = b + \sum_i x_i * w_i$$



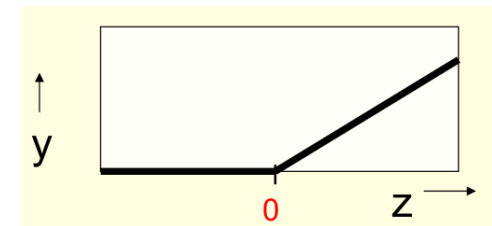
Artificial neuron activations

Then apply activation function:

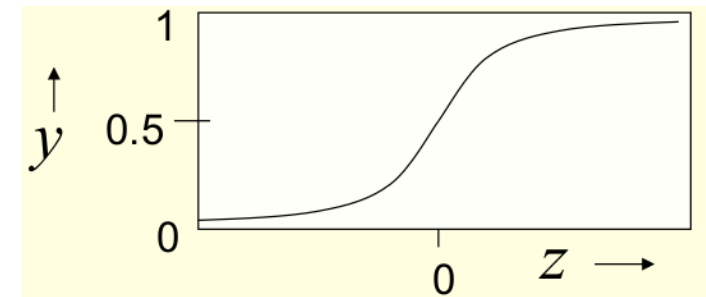
- Linear neuron: $y = z$
- Binary threshold neuron: $y = \begin{cases} 1, z \geq 0 \\ 0, z < 0 \end{cases}$



- Rectified Linear neuron: $y = \max(0, z)$



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



- Tanh neuron: $y = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

A machine learning algorithm
usually corresponds to a combination of
the following 3 elements:
(either explicitly specified or implicit)

- ✓ the choice of a specific **function family**: F
(often a parameterized family)
- ✓ a **way to evaluate the quality** of a function $f \in F$
(typically using a **cost (or loss) function** L
measuring how wrongly f predicts)
- ✓ a **way to search for the «best»** function $f \in F$
(typically an optimization of function parameters to
minimize the overall loss over the training set).

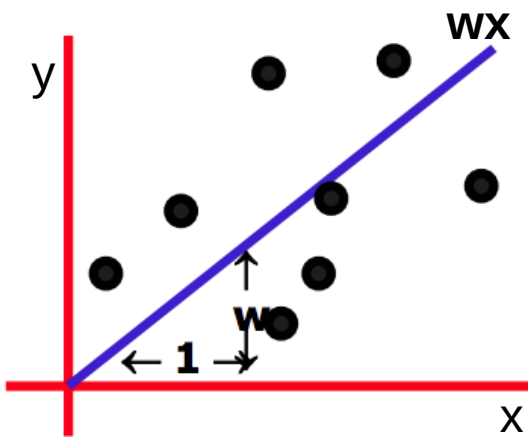
1D Linear Regression

Regression: for input vector x predict real value y (cf. Classification)

- Approximate dependence with linear function:

$$\hat{y} = h_w(x) = wx \rightarrow \text{Hypothesis}$$

- Learn w from data: Loss function? Optimization algorithm?



inputs	outputs
$x_1 = 1$	$y_1 = 1$
$x_2 = 3$	$y_2 = 2.2$
$x_3 = 2$	$y_3 = 2$
$x_4 = 1.5$	$y_4 = 1.9$
$x_5 = 4$	$y_5 = 3.1$

1D / Simple Regression:
Regression with 1 input variable (feature)

1D Linear Regression Loss

Loss (Cost, Objective) function:

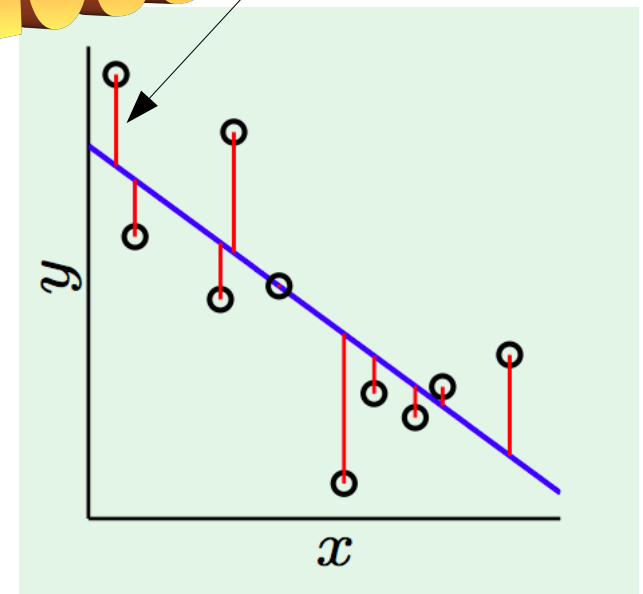
- how well hypothesis $h(x)$ approximates real function $y=f(x)$
- minimize loss on train set and hope it will generalize

For regression use sum of squared **residuals/errors**:

$$E = \sum_i (y_i - wx_i)^2$$

Loss

- Natural
- Justified by Maximum Likelihood principle under normal distribution of residuals ← **DERIVE!**



1D Linear Regression Optimization

E can be minimized analytically:

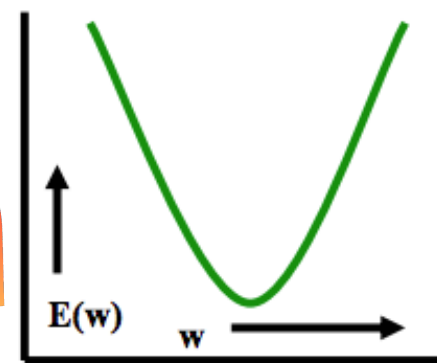
- not the case for more sophisticated models

1. $E(w)$ is convex \Rightarrow single optimum=global minimum

2. $E'(w) = 0 \leftarrow$ **SOLVE**

$$w = \frac{\sum x_i y_i}{\sum x_i^2}$$

Optimization



$$\begin{aligned} E &= \sum_i (y_i - wx_i)^2 \\ &= \sum_i y_i^2 - (2 \sum_i x_i y_i)w + (\sum_i x_i^2)w^2 \end{aligned}$$

Multivariate Linear Regression

$$\hat{y} = h_w(x) = w^T x = \sum_{j=1}^m w_j * x_j \quad \rightarrow$$

Hypothesis

Matrix-vector product: much faster than loop! **always VECTORIZED computations**

Loss

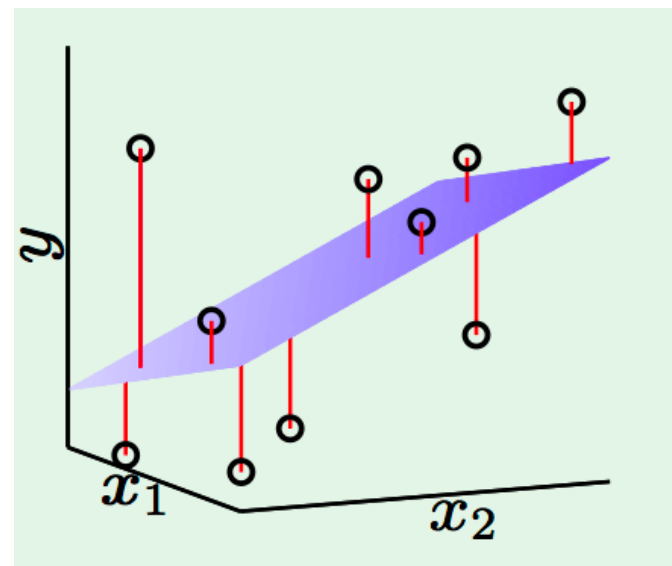
Mean squared error (MSE):

Multiplying by constant doesn't change argmin, but easier to compare

$$E_{\text{in}}(w) = \frac{1}{N} \sum_{i=1}^N (w^T x_{\{i\}} - y_{\{i\}})^2$$

$$E_{\text{in}}(w) = \frac{1}{N} \|Xw - y\|^2$$

$$X = \begin{bmatrix} -\mathbf{x}_1^T- \\ -\mathbf{x}_2^T- \\ \vdots \\ -\mathbf{x}_N^T- \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$



Multivariate Linear Regression Optimization

Loss is convex, single optimum=global minimum

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{2}{N} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^\dagger \mathbf{y} \quad \text{where} \quad \mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

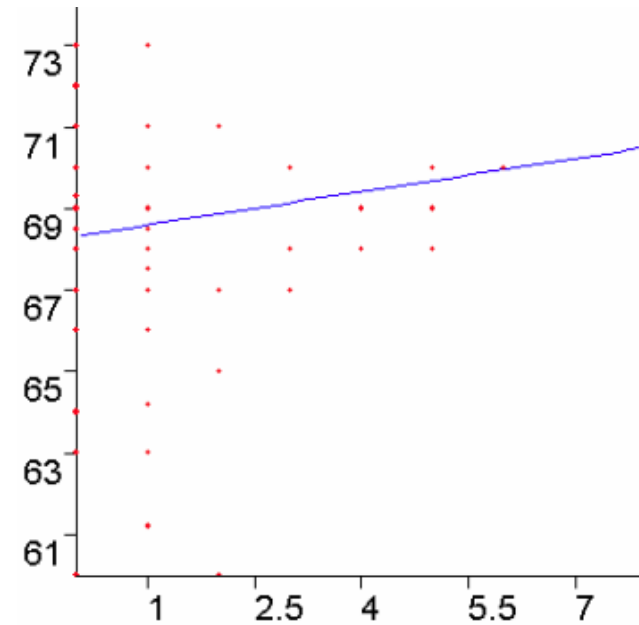
\mathbf{X}^\dagger is the 'pseudo-inverse' of \mathbf{X}

Optimization

Bias / intercept

What if $f(0) \neq 0$?

- Add bias / intercept: $h(x) = wx + b$
- Concatenate “always 1” input to simplify implementation

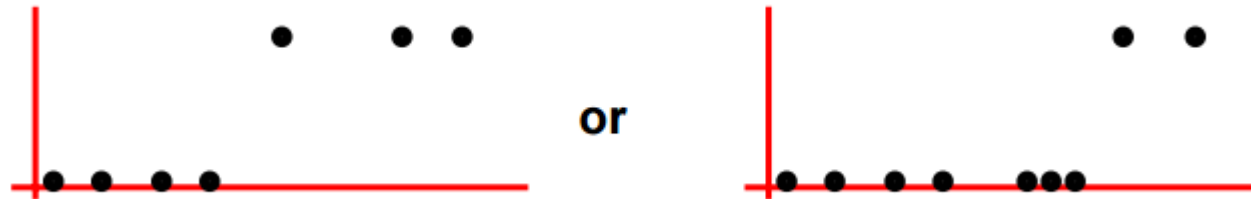


$$\hat{y} = h_w(x) = b + \sum_{j=1}^m w_j * x_j = [1; x]^T w, \text{ where } w_0 = b$$

x_0	x_1	x_2	y
1	2	4	16
1	3	4	17
1	5	5	20

Classification with linear regression

Classification as regression where y is 0 or 1



- Train linear regression

- Predict
$$\begin{cases} 1, \hat{y} \geq 0.5 \\ 0, \hat{y} < 0.5 \end{cases}$$



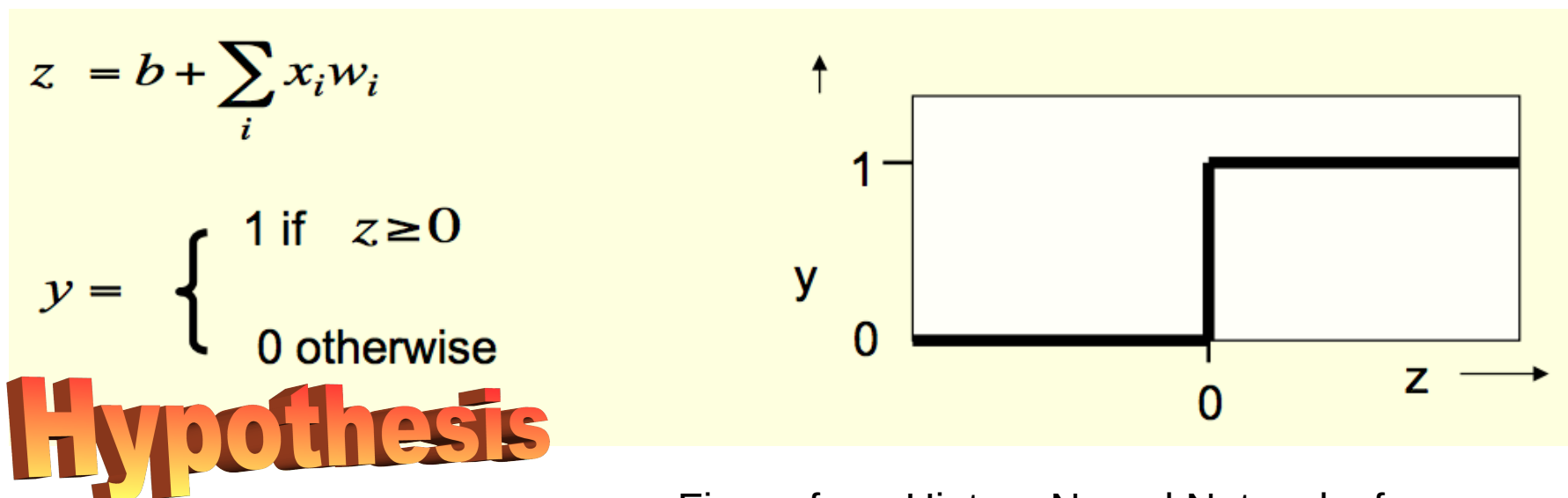
Very unstable decision boundary!

- Least squares tries hard to minimize residuals for points which are already classified correctly!

Classification with perceptron

Binary threshold neuron – early mathematical model a biological neuron

- Proposed by McCulloch & Pitts in 1943
- Biological neuron output was supposed to represent truth of some logical proposition



Classification with perceptron

In 1960s Frank Rosenblatt proposed a learning algorithm for perceptron:

Use -1 / +1 instead of 0 / 1 as labels

until convergence:

- pick next (x_i, y_i)
- If correctly classified: do nothing
- If incorrectly classified as 1: $w := w - x$
- If incorrectly classified as -1: $w := w + x$



Pros: If dataset is linearly separable, this algorithm is guaranteed to converge to some solution ← PROOF in Hinton's lecture2@coursera

Cons: very slow convergence, no guarantees for non linearly separable data

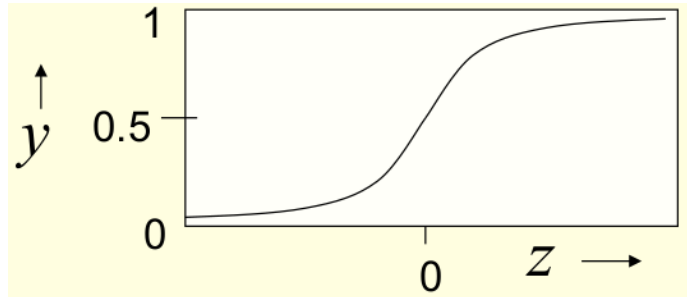
It was (incorrectly) shown that perceptrons can classify photos of tanks vs. trucks obscured in the forest

Photos of trucks were taken on a sunny day, photos of tanks – in a cloudy day, perceptron simply learned to estimate brightness.

Logistic regression (LR)

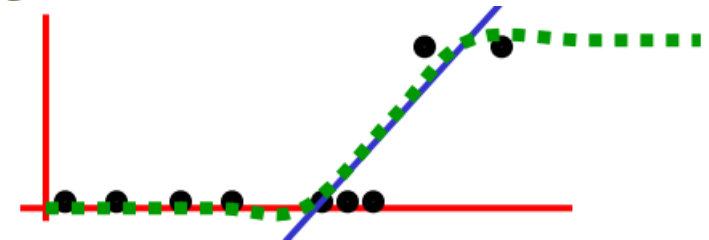
- Use sigmoid neuron instead of linear

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



$$\hat{y} = h_w(x) = \sigma([1; x]^T w) \longrightarrow \text{Hypothesis}$$

- more stable decision boundary
- Outputs value in (0,1) – can be interpreted as probability of a positive class



Logistic regression (LR) Loss

- Use binary cross-entropy loss (CE)
 - LR+CE is convex w.r.t. weights (unlike LR+MSE)
 - Justified by Maximum Likelihood ← DERIVE

Loss

$$E(w) = -\frac{1}{N} \sum_{i=1}^N y_{\{i\}} \log(h_w(x_{\{i\}})) + (1 - y_{\{i\}}) \log(1 - h_w(x_{\{i\}}))$$

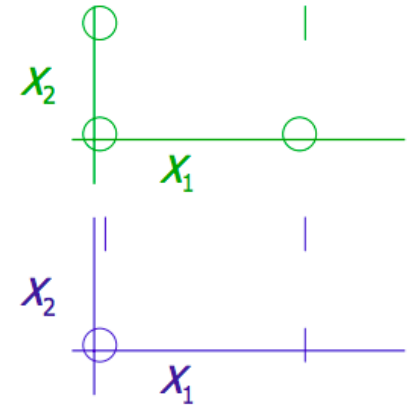
- The same loss is commonly used for classification using NNs!

XOR problem $x_1 \wedge x_2$

Linear models can learn:

- **AND, OR** functions \leftarrow FIND W ?
- Any conjunction or disjunction of literals and their negations
 - x_1 AND NOT x_2 AND x_3
 - x_1 OR NOT x_2 OR x_3

$$x_1 \vee x_2$$



Linear models cannot learn **XOR** function!

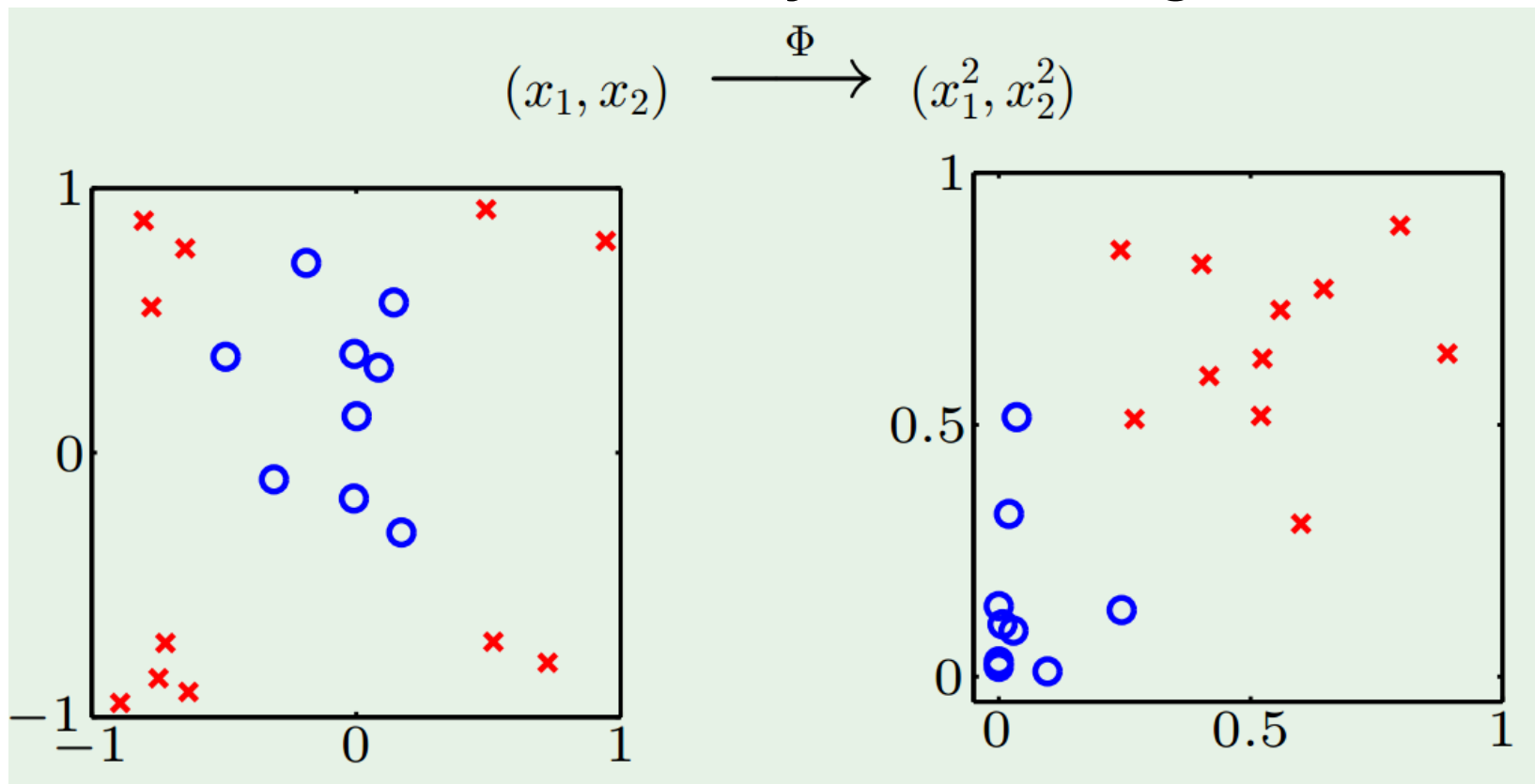
- Linearly non-separable
- People incorrectly concluded that NNs are not suited for real tasks
- NN with 1 hidden layer can solve XOR \leftarrow SHOW

Modeling non-linear relations

Linear models can be used for non-linear relations

- but needs hand feature engineering

Linear models = linearity w.r.t. weights, not features

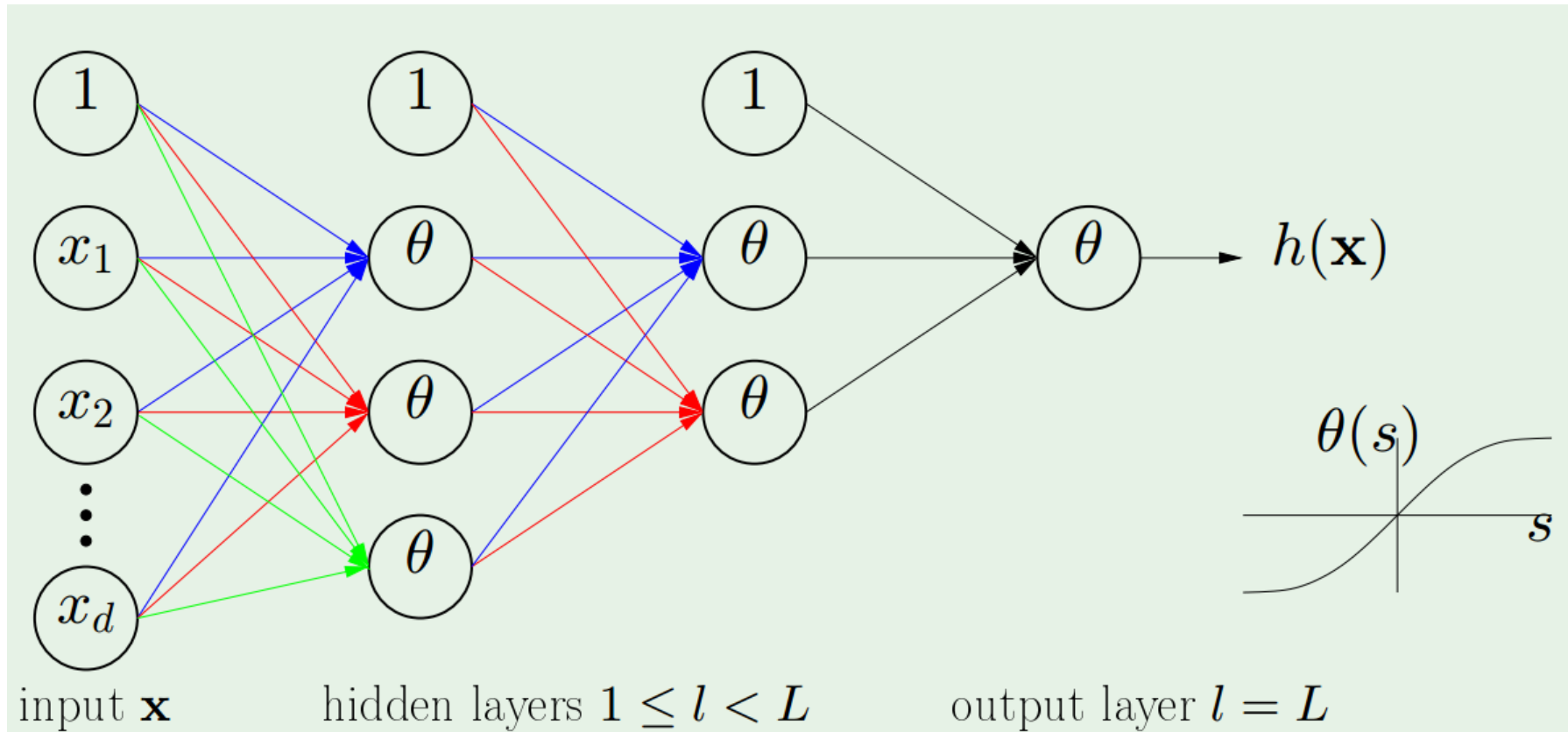


Feed-forward NN (FFNN)

FFNN – simply a composition of logistic regressions!

Even 1 hidden layer FFNN (given enough hidden units) can represent:

- any boolean function (exactly)
 - ← BUILD 1 hidden layer NN with binary threshold neurons for XOR (manually)!
- any continuous function defined on compact subsets of \mathbb{R}^n (approximately with any precision)
Appropriate weights exist, but no guarantees they will be learned!

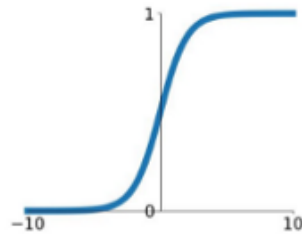


Popular activation functions

- Don't use Sigmoid FFNN (not centered in 0)!
- Tanh is good default

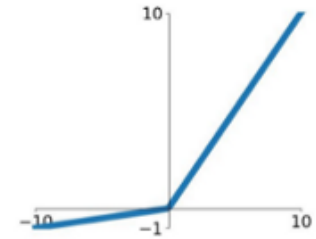
Sigmoid

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



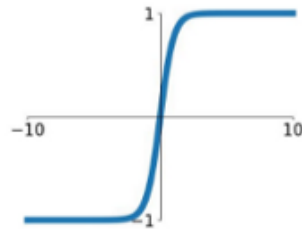
Leaky ReLU

$$\max(0.1x, x)$$



tanh

$$\tanh(x)$$

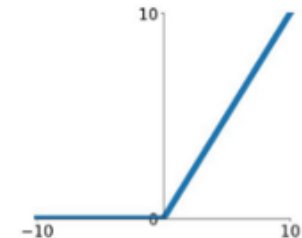


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

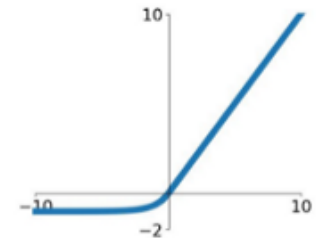
ReLU

$$\max(0, x)$$



ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



LR/FFNN Optimization

- Optimized using iterative optimization algorithms (not analytically)
 - Most commonly used is Stochastic Gradient Descent (SGD) ← next lecture

