

SD

Mechatronic Engineering program

Basics of AI and Deep Learning:
4: Artificial Neural Networks
 – in three flavors

Ziemowit Dworakowski
 AGH University of Krakow

Understanding is measured by how many different perspectives you can internalize for a given subject

1

SD

So far, we've learned:

- What is classification and regression – from the data perspective
- How to use optimization methods for classifier or regressor training
- What is a feature space and what can we find in it...

Today, we'll learn about neural networks

Understanding is measured by how many different perspectives you can internalize for a given subject...

2

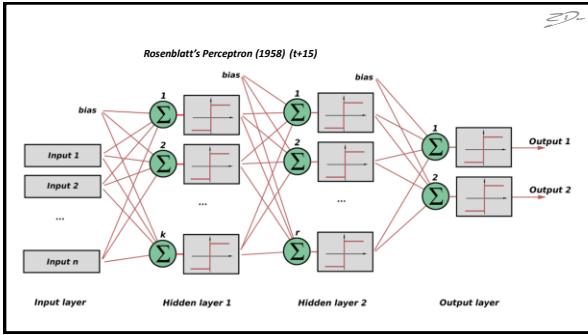
SD

Lets model a neuron:

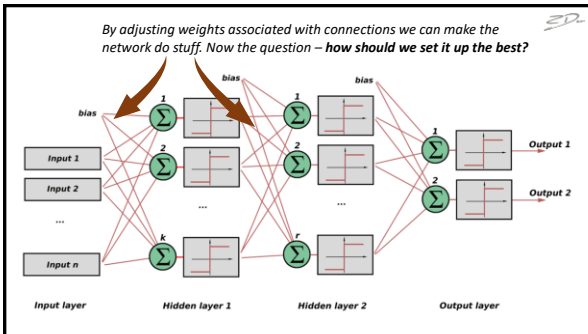
McCulloch – Pitts neuron (1943)

One neuron is not really that brilliant... But maybe we can connect it with others?

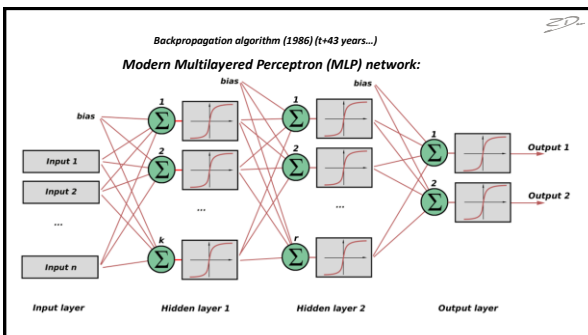
3



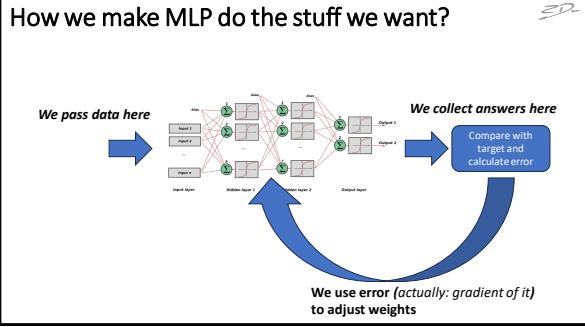
4



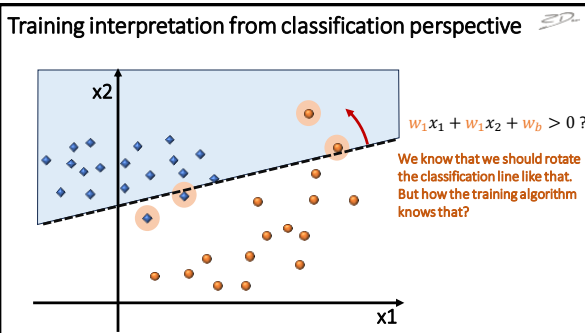
5



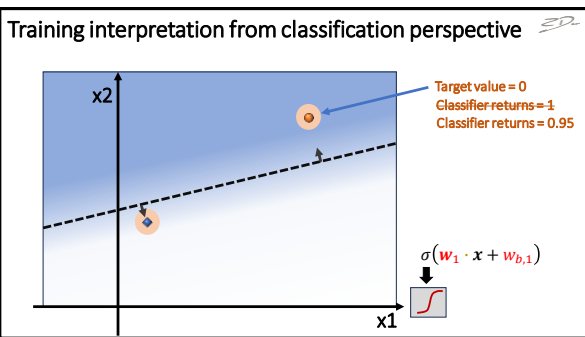
6



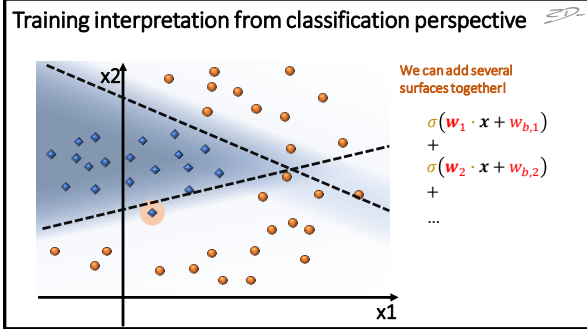
7



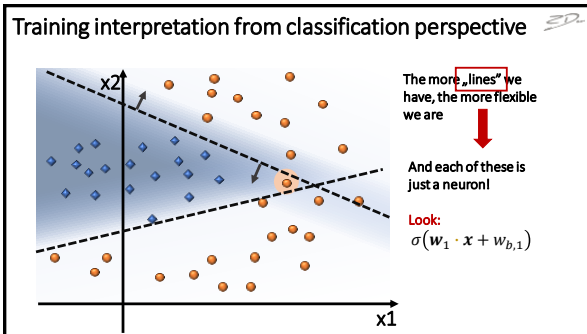
8



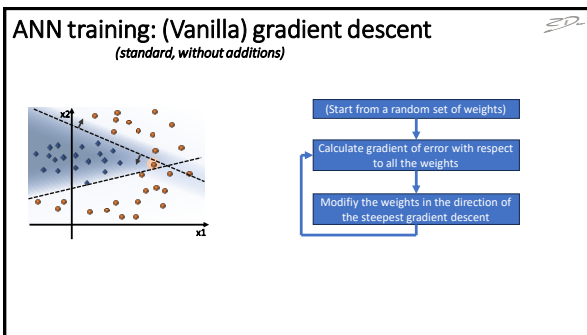
9



10



11



12

Data space mapping function SD

$y = f(\mathbf{w}, \mathbf{x})$

*We'll need lots of flexibility and lots of non-linearity...
And we want it to be adjustable according to needs...*

13

$y = f(\mathbf{w}, \mathbf{x})$ SD

We might want to make it compounding from different functions:
 $y = f_1(\mathbf{w}_1, \mathbf{x}) + f_2(\mathbf{w}_2, \mathbf{x}) + \dots + f_n(\mathbf{w}_n, \mathbf{x})$

We might want to govern these functions purely by their parameters:
 $x'_1 = \sigma(\mathbf{w}_{1,1}, \mathbf{x}) + \sigma(\mathbf{w}_{2,1}, \mathbf{x}) + \dots + \sigma(\mathbf{w}_{n,1}, \mathbf{x})$
 $x'_2 = \sigma(\mathbf{w}_{1,2}, \mathbf{x}) + \sigma(\mathbf{w}_{2,2}, \mathbf{x}) + \dots + \sigma(\mathbf{w}_{n,2}, \mathbf{x})$
 \dots

← *This is a vector of our new coordinates!*

And now a „final decision“ – another nonlinear transformation:
 $y = \sigma(\mathbf{w}_{n+1,1}, \mathbf{x}')$

If we want, we can nest it even further (so y becomes x'' and so on)

14

Data space mapping function SD

$y = f(\mathbf{w}, \mathbf{x})$

*And this tool for transformation of feature space (x) into
output space (y) we'll call a **neural network**:*
 $y = \sigma(\mathbf{w}_{n+1,1}, \mathbf{x}')$

We adjust it to data using a gradient-descent backpropagation algorithm

15

Three flavors of neural networks:

A **compound function** that does complex non-linear mapping from input space into output space morphing data to gradually simplify its structure

A set of **hyperplanes** in feature space (whose location we can adjust) that form layered decision surfaces

A **biologically-inspired network of nodes** that can pass information through one another with **emergent feature** of data understanding

neural network is:

16

Consider a simple Multilayer Perceptron (MLP) network:

We also have a target t associated with each input vector x . Using that, we can calculate error d :

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

Now we can calculate gradient of d with respect to weights w :

$$\nabla d = \left[\frac{\partial d}{\partial w_1}, \frac{\partial d}{\partial w_2}, \frac{\partial d}{\partial w_3}, \dots, \frac{\partial d}{\partial w_9} \right]$$

The gradient now tells us how to adjust weights

$y = \sigma(w_8 \cdot x'_1 + w_9 \cdot x'_2 + w_7)$

$x'_1 = \sigma(w_3 \cdot x_1 + w_5 \cdot x_2 + w_1)$

$x'_2 = \sigma(w_2 \cdot x_1 + w_6 \cdot x_2 + w_2)$

You can combine this to get a $f(w, x)$ function equivalent for a network

17

Consider a simple Multilayer Perceptron (MLP) network:

We also have a target t associated with each input vector x . Using that, we can calculate error d :

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

Now we can calculate gradient of d with respect to weights w :

$$\nabla d = \left[\frac{\partial d}{\partial w_1}, \frac{\partial d}{\partial w_2}, \frac{\partial d}{\partial w_3}, \dots, \frac{\partial d}{\partial w_9} \right]$$

The gradient now tells us how to adjust weights

$y = \sigma(w_8 \cdot x'_1 + w_9 \cdot x'_2 + w_7)$

$\frac{\partial d}{\partial w_9} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} \frac{\partial p}{\partial w_9} = (y - t)y(1 - y) \cdot x'_2$

$\frac{\partial d}{\partial y} = y - t$

18

Consider a simple Multilayer Perceptron (MLP) network: *We also have a target t associated with each input vector x . Using that, we can calculate error d :*

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

Now we can calculate gradient of d with respect to weights w :

$$\nabla d = \left[\frac{\partial d}{\partial w_1}, \frac{\partial d}{\partial w_2}, \frac{\partial d}{\partial w_3}, \dots, \frac{\partial d}{\partial w_9} \right]$$

The gradient now tells us how to adjust weights

$y = \sigma(w_8 \cdot x'_1 + w_9 \cdot x'_2 + w_7)$

$$\frac{\partial d}{\partial w_9} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} \frac{\partial p}{\partial w_9} = (y - t)y(1 - y) \cdot x'_2$$

$y(1 - y) \leftarrow y = \sigma(p) = \frac{1}{1 + e^{-p}}$

19

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

$$\frac{\partial \sigma}{\partial p} = \frac{\partial}{\partial p} (1 + e^{-p})^{-1} = -1(1 + e^{-p})^{-2} \cdot e^{-p} \cdot -1$$

$$\frac{\partial \sigma}{\partial p} = \frac{1}{1 + e^{-p}} \cdot \frac{1}{1 + e^{-p}} \cdot e^{-p}$$

$$\frac{\partial \sigma}{\partial p} = \frac{e^{-p}}{1 + e^{-p}} \cdot \frac{1}{1 + e^{-p}}$$

$$\frac{\partial \sigma}{\partial p} = \frac{1 + e^{-p} - 1}{1 + e^{-p}} \cdot \frac{1}{1 + e^{-p}}$$

$$\frac{\partial \sigma}{\partial p} = \left(1 - \frac{1}{1 + e^{-p}} \right) \cdot \frac{1}{1 + e^{-p}} = (1 - y) \cdot y$$

20

Consider a simple Multilayer Perceptron (MLP) network: *We also have a target t associated with each input vector x . Using that, we can calculate error d :*

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

Now we can calculate gradient of d with respect to weights w :

$$\nabla d = \left[\frac{\partial d}{\partial w_1}, \frac{\partial d}{\partial w_2}, \frac{\partial d}{\partial w_3}, \dots, \frac{\partial d}{\partial w_9} \right]$$

The gradient now tells us how to adjust weights

$y = \sigma(w_8 \cdot x'_1 + w_9 \cdot x'_2 + w_7)$

$$\frac{\partial d}{\partial w_9} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} \frac{\partial p}{\partial w_9} = (y - t)y(1 - y) \cdot x'_2$$

$\frac{\partial p}{\partial w_9} = x'_2$

21

Consider a simple Multilayer Perceptron (MLP) network:

We also have a target t associated with each input vector x . Using that, we can calculate error d :

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

Now we can calculate gradient of d with respect to weights w :

$$\nabla d = \left[\frac{\partial d}{\partial w_1}, \frac{\partial d}{\partial w_2}, \frac{\partial d}{\partial w_3}, \dots, \frac{\partial d}{\partial w_9} \right]$$

The gradient now tells us how to adjust weights

$$\frac{\partial d}{\partial w_9} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} \frac{\partial p}{\partial w_9} = (y - t)y(1 - y) \cdot x'_2$$

$$\frac{\partial p}{\partial w_9} = x'_2$$

22

Consider a simple Multilayer Perceptron (MLP) network:

We also have a target t associated with each input vector x . Using that, we can calculate error d :

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

Now we can calculate gradient of d with respect to weights w :

$$\nabla d = \left[\frac{\partial d}{\partial w_1}, \frac{\partial d}{\partial w_2}, \frac{\partial d}{\partial w_3}, \dots, \frac{\partial d}{\partial w_9} \right]$$

The gradient now tells us how to adjust weights

As we go through the network backwards, we propagate our error calculations – it is the **backpropagation algorithm**

$$\frac{\partial d}{\partial w_9} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} \frac{\partial p}{\partial w_9} = (y - t)y(1 - y) \cdot x'_2$$

$$\frac{\partial d}{\partial w_6} = \frac{\partial d}{\partial y} \frac{\partial y}{\partial p} \frac{\partial p}{\partial x'_2} \frac{\partial x'_2}{\partial w_6} = (y - t)y(1 - y) \cdot x'_2 \cdot u$$

23

We stop training when:

- Validation error starts to rise
- We run out of time
- Gradient approaches 0 (we don't improve any more)

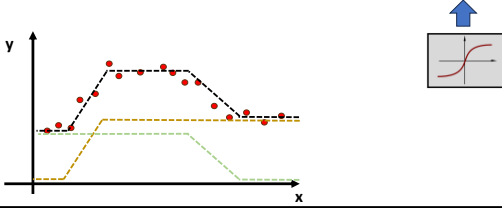
Epoch means 1 passing of the entire training dataset through training

24

Regression with multilayered perceptron

The idea here is to take a lot of simple nonlinear functions and add them together to get a more complex one

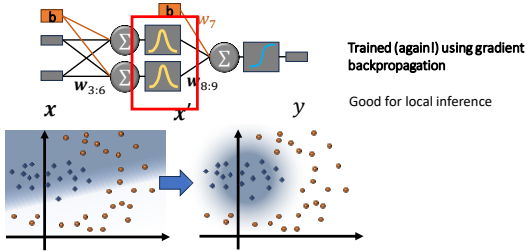
$$f(\mathbf{w}, \mathbf{x}) = f_1(\mathbf{w}_1, \mathbf{x}) + f_2(\mathbf{w}_2, \mathbf{x}) + \dots + f_3(\mathbf{w}_3, \mathbf{x})$$



25

Radial networks

SD

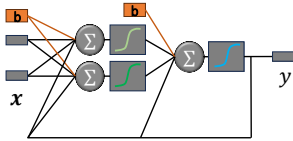


Trained (again!) using gradient backpropagation
Good for local inference

26

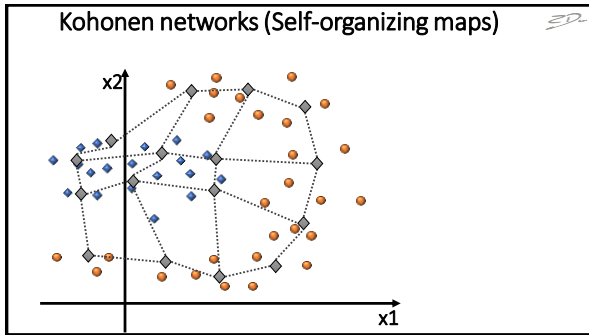
Recurrent networks

SD

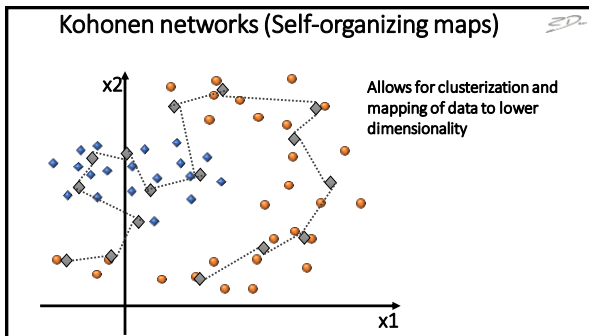


Adding recurrent connections makes for good short-time memory capabilities – allowing for system identification and prediction of time-domain signals

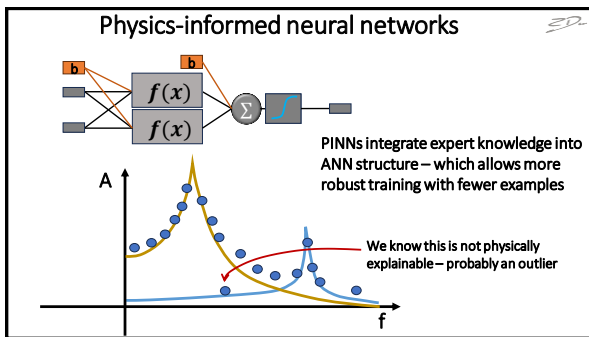
27



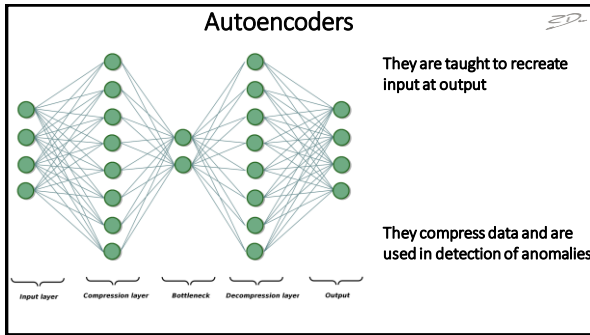
28



29



30



31

- ### Things to remember: SD
1. Explain Multilayered Perceptron Networks using one of the three explanations
 2. Draw and annotate McCulloch-Pitts neuron model
 3. Explain why McCulloch-Pitts neuron is represented by a straight line in a feature space
 4. Draw and annotate full scheme of a Multilayer Perceptron Network
 5. Explain the idea of a gradient backpropagation algorithm (no need to memorize calculations)
 6. Write equation for a sigmoid activation function
 7. Explain basic procedure that allows to prevent MLP network from overfitting (how do we use training and validation data to this end?)
 8. Explain general idea behind selected network structures different than MLP (Kohonen, PINN, Radial, Autoencoders, Recurrent)

32
