1/40

#### Moderne Methoden der Datenanalyse

#### Depth and activation?

Roger Wolf 30. June 2023

# **Content of this lecture**

- Recap: perceptron  $\rightarrow$  multilayer perceptron (MLP).
- MLP as representation of Boolean functions.
- MLP as representation of arbitrary contours.
- MLP training:
  - Perceptron learning rule.
  - Not linearly separable tasks and activation functions revisted.

#### <sup>3/40</sup> Recap: perceptron

• Last time we discussed the Boolean perceptron with a simple step function as **activation function** (here shown for real-valued *input features*):



### <sup>4/40</sup> **Recap: activation function**

- We discussed that the activation function could be any function.
- A few popular examples are given below:

#### Rectified linear unit (ReLU):

$$\theta(z) = \begin{cases} z & z \ge 0\\ 0 & \text{sonst} \end{cases}$$

#### Sigmoid:



#### tanh:

 $\theta(z) = \tanh(z)$ 



#### Softplus:

 $\theta(z) = \log(1 + \exp(z))$ 



### <sup>5/40</sup> Recap: multilayer perceptron (MLP)

 We have discussed that the ability to express Boolean relations increases when using more than one perceptron, organized in layers → multilayer perceptron (MLP):



### The fully-connected feed-forward NN

- Fully-connected: All nodes of consecutive layers are *connected* with each other.
- Feed-forward: Inputs are propagated only in *forward* direction.



#### 7/40 Mathematical notation



#### <sup>8/40</sup> **Depth**

- A feed-forward NN can be understood as a **directed graph** of depth *d*.
- A directed graph has *sources* and *drains*. The depth of a graph is the longest path between a source and a drain.



• An NN with a depth of d > 2 (i.e. an NN with more than 2 hidden layers) we call *deep*.

# <sup>9/40</sup> The MLP as representation of Boolean functions

• Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality?

## <sup>9/40</sup> The MLP as representation of Boolean functions

 Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality? – 1.

- Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality? – 1.
- **Proof**: any Boolean function can be expressed in the form of a **truth table**.

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | y |
|-------|-------|-------|-------|-------|---|
| 0     | 0     | 1     | 1     | 0     | 1 |
| 0     | 1     | 0     | 1     | 1     | 1 |
| 0     | 1     | 1     | 0     | 0     | 1 |
| 1     | 0     | 0     | 0     | 1     | 1 |
| 1     | 0     | 1     | 1     | 1     | 1 |
| 1     | 1     | 0     | 0     | 1     | 1 |

Arbitrary example

- Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality? – 1.
- **Proof**: any Boolean function can be expressed in the form of a **truth table**.

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | <i>y</i> — |
|-------|-------|-------|-------|-------|------------|
| 0     | 0     | 1     | 1     | 0     | 1          |
| 0     | 1     | 0     | 1     | 1     | 1          |
| 0     | 1     | 1     | 0     | 0     | 1          |
| 1     | 0     | 0     | 0     | 1     | 1          |
| 1     | 0     | 1     | 1     | 1     | 1          |
| 1     | 1     | 0     | 0     | 1     | 1          |

Disjunctive normal form (DNF):

 $y = \overline{x}_1 \overline{x}_2 x_3 x_4 \overline{x}_5 \, \lor \, \overline{x}_1 x_2 \overline{x}_3 x_4 x_5 \, \lor \overline{x}_1 x_2 x_3 \overline{x}_4 \overline{x}_5 \, \lor$ 

 $x_1\overline{x}_2\overline{x}_3\overline{x}_4x_5 \lor x_1\overline{x}_2x_3x_4x_5 \lor x_1x_2\overline{x}_3\overline{x}_4x_5$ 

Arbitrary example

- Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. • How many hidden layers does the MLP minimally require to have this quality? – 1.
- **Proof**: any Boolean function can be • expressed in the form of a truth table.

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | y |
|-------|-------|-------|-------|-------|---|
| 0     | 0     | 1     | 1     | 0     | 1 |
| 0     | 1     | 0     | 1     | 1     | 1 |
| 0     | 1     | 1     | 0     | 0     | 1 |
| 1     | 0     | 0     | 0     | 1     | 1 |
| 1     | 0     | 1     | 1     | 1     | 1 |
| 1     | 1     | 0     | 0     | 1     | 1 |





 $x_3$ 

 $x_5$ 

 $x_2$ 

 $x_1$ 

- Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality? – 1.
- **Proof**: any Boolean function can be expressed in the form of a **truth table**.

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | y |
|-------|-------|-------|-------|-------|---|
| 0     | 0     | 1     | 1     | 0     | 1 |
| 0     | 1     | 0     | 1     | 1     | 1 |
| 0     | 1     | 1     | 0     | 0     | 1 |
| 1     | 0     | 0     | 0     | 1     | 1 |
| 1     | 0     | 1     | 1     | 1     | 1 |
| 1     | 1     | 0     | 0     | 1     | 1 |





- Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality? – 1.
- **Proof**: any Boolean function can be expressed in the form of a **truth table**.

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | y |
|-------|-------|-------|-------|-------|---|
| 0     | 0     | 1     | 1     | 0     | 1 |
| 0     | 1     | 0     | 1     | 1     | 1 |
| 0     | 1     | 1     | 0     | 0     | 1 |
| 1     | 0     | 0     | 0     | 1     | 1 |
| 1     | 0     | 1     | 1     | 1     | 1 |
| 1     | 1     | 0     | 0     | 1     | 1 |

Arbitrary example



- Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality? – 1.
- **Proof**: any Boolean function can be expressed in the form of a **truth table**.

| $x_1$     | $x_2$   | $x_3$ | $x_4$ | $x_5$ | y |   |
|-----------|---------|-------|-------|-------|---|---|
| 0         | 0       | 1     | 1     | 0     | 1 | $y = \underbrace{\overline{x}_1 \overline{x}_2 x_3 x_4 \overline{x}}_{(1)}$ |
| 0         | 1       | 0     | 1     | 1     | 1 | $\bigcup_{x_1\overline{x}_2\overline{x}_3\overline{x}_4x}$                  |
| 0         | 1       | 1     | 0     | 0     | 1 | 4   |
| 1         | 0       | 0     | 0     | 1     | 1 |   |
| 1         | 0       | 1     | 1     | 1     | 1 |   |
| 1         | 1       | 0     | 0     | 1     | 1 |   |
| Arbitrary | example |       |       |       |   | $\backslash$  |



- Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality? – 1.
- **Proof**: any Boolean function can be expressed in the form of a **truth table**.

| $x_1$     | $x_2$   | $x_3$ | $x_4$ | $x_5$ | y |  |
|-----------|---------|-------|-------|-------|---|--|
| 0         | 0       | 1     | 1     | 0     | 1 | $y = \underbrace{\overline{x_1}\overline{x_2}x_3x_4\overline{x_5}}_{(1)} \lor \underbrace{\overline{x_1}x_2\overline{x_3}x_4x_5}_{(2)} \lor \underbrace{\overline{x_1}x_2x_3\overline{x_4}\overline{x_5}}_{(2)} \lor$  |
| 0         | 1       | 0     | 1     | 1     | 1 | $\underbrace{(1)}_{x_1\overline{x}_2\overline{x}_3\overline{x}_4x_5} \vee \underbrace{(x_1\overline{x}_2x_3x_4x_5)}_{x_1\overline{x}_2x_3\overline{x}_4x_5} \vee \underbrace{(x_1\overline{x}_2\overline{x}_3\overline{x}_4x_5)}_{x_1\overline{x}_2\overline{x}_3\overline{x}_4x_5}$ |
| 0         | 1       | 1     | 0     | 0     | 1 | <u>(4)</u> <u>(5)</u>  |
| 1         | 0       | 0     | 0     | 1     | 1 |  |
| 1         | 0       | 1     | 1     | 1     | 1 |  |
| 1         | 1       | 0     | 0     | 1     | 1 |  |
| Arbitrary | example |       |       |       |   |  |
|           |         |       |       |       |   | $x_1  x_2  x_3  x_4  x_5$  |

- Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality? – 1.
- **Proof**: any Boolean function can be expressed in the form of a **truth table**.

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | y |
|-------|-------|-------|-------|-------|---|
| 0     | 0     | 1     | 1     | 0     | 1 |
| 0     | 1     | 0     | 1     | 1     | 1 |
| 0     | 1     | 1     | 0     | 0     | 1 |
| 1     | 0     | 0     | 0     | 1     | 1 |
| 1     | 0     | 1     | 1     | 1     | 1 |
| 1     | 1     | 0     | 0     | 1     | 1 |

Arbitrary example



- Last time we have discussed that an MLP can represent <u>any arbitrary Boolean function</u>. How many hidden layers does the MLP minimally require to have this quality? – 1.
- **Proof**: any Boolean function can be expressed in the form of a **truth table**.

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | y |
|-------|-------|-------|-------|-------|---|
| 0     | 0     | 1     | 1     | 0     | 1 |
| 0     | 1     | 0     | 1     | 1     | 1 |
| 0     | 1     | 1     | 0     | 0     | 1 |
| 1     | 0     | 0     | 0     | 1     | 1 |
| 1     | 0     | 1     | 1     | 1     | 1 |
| 1     | 1     | 0     | 0     | 1     | 1 |

Arbitrary example



**[Without proof]** For N inputs the number of required nodes can grow exponentially (exp(N)). Giving the NN more depth, it is possible to reduce the number of nodes down to  $log_2(N)$ !

| $x_1$     | $x_2$   | $x_3$ | $x_4$ | $x_5$ | y |   |
|-----------|---------|-------|-------|-------|---|---|
| 0         | 0       | 1     | 1     | 0     | 1 | $y = \underbrace{\overline{x_1}\overline{x_2}x_3x_4\overline{x_5}}_{(1)} \lor \underbrace{\overline{x_1}x_2\overline{x_3}x_4x_5}_{(2)} \lor \underbrace{\overline{x_1}x_2x_3\overline{x_4}\overline{x_5}}_{(2)} \lor$   |
| 0         | 1       | 0     | 1     | 1     | 1 | $\underbrace{(1)}_{x_1\overline{x}_2\overline{x}_3\overline{x}_4x_5} \vee \underbrace{(x_1\overline{x}_2x_3x_4x_5)}_{x_1\overline{x}_2x_3\overline{x}_4x_5} \vee \underbrace{(x_1x_2\overline{x}_3\overline{x}_4x_5)}_{x_1\overline{x}_2\overline{x}_3\overline{x}_4x_5}$ |
| 0         | 1       | 1     | 0     | 0     | 1 | $\underbrace{\underbrace{4}}_{5}$   |
| 1         | 0       | 0     | 0     | 1     | 1 |   |
| 1         | 0       | 1     | 1     | 1     | 1 |   |
| 1         | 1       | 0     | 0     | 1     | 1 |   |
| Arbitrary | example |       |       |       |   | $\begin{array}{c ccccccccccccccccccccccccccccccccccc$   |

Priv.-Doz. Dr. Roger Wolf http://ekpwww.physik.uni-karlsruhe.de/~rwolf/

• We have discussed that an MLP can approximate any *arbitrary boundary* (like the one shown below) with *abitrary precision*.



- For this we have used two hidden layers.
- Question: Do you think that this can be done in general<sup>[1]</sup> with only one single hidden layer? If yes how? If no why not?

• Answer: it is possible, but a general construction is more complicated as you might think.

- Answer: it is possible, but a general construction is more complicated as you might think.
- The root cause is the difficulty to represent an arbitrary convex boundary with an MLP:
- A contour with N=4 bounds can be represented by an MLP with a single hidden layer as shown (for a single pentagon) e.g. on slide 15.





- Answer: it is possible, but a general construction is more complicated as you might think.
- The root cause is the difficulty to represent an arbitrary convex boundary with an MLP:
- A contour with N=4 bounds can be represented by an MLP with a single hidden layer as shown (for a single pentagon) e.g. on slide 15.
- What prevents us from filling an arbitrary contour with small squares? there is an unbound area with MLP output 0 < y < N, i.e. when filling the boundary with squares a distinction between "in- and outside the contour" is impossible!</li>



MLP with 4 nodes

- Answer: it is possible, but a general construction is more complicated as you might think.
- The root cause is the difficulty to represent an arbitrary convex boundary with an MLP:
- A contour with N=4 bounds can be represented by an MLP with a single hidden layer as shown (for a single pentagon) e.g. on slide 15.
- What prevents us from filling an arbitrary contour with small squares? there is an unbound area with MLP output 0 < y < N, i.e. when filling the boundary with squares a distinction between "in- and outside the contour" is impossible!</li>
- This issue can be mitigated by moving on to more bounds.



- Answer: it is possible, but a general construction is more complicated as you might think.
- The root cause is the difficulty to represent an arbitrary convex boundary with an MLP:
- A contour with N=5 bounds can be represented by an MLP with a single hidden layer as shown (for a single pentagon) e.g. on slide 15.



- Answer: it is possible, but a general construction is more complicated as you might think.
- The root cause is the difficulty to represent an arbitrary convex boundary with an MLP:
- A contour with N=6 bounds can be represented by an MLP with a single hidden layer as shown (for a single pentagon) e.g. on slide 15.



- Answer: it is possible, but a general construction is more complicated as you might think.
- The root cause is the difficulty to represent an arbitrary convex boundary with an MLP:
- A contour with N=8 bounds can be represented by an MLP with a single hidden layer as shown (for a single pentagon) e.g. on slide 15.



- Answer: it is possible, but a general construction is more complicated as you might think.
- The root cause is the difficulty to represent an arbitrary convex boundary with an MLP:
- A contour with N=8 bounds can be represented by an MLP with a single hidden layer as shown (for a single pentagon) e.g. on slide 15.





- Answer: it is possible, but a general construction is more complicated as you might think.
- In this way any arbitrary contour can be approximated with arbitrary precision, with an MLP with a single hidden layer by filling the contours with cylinders. The perceptrons of the hidden layer are added in the output layer as demonstrated, e.g., on slide 15:



- Answer: it is possible, but a general construction is more complicated as you might think.
- In this way any arbitrary contour can be approximated with arbitrary precision, with an MLP with a single hidden layer by filling the contours with cylinders. The perceptrons of the hidden layer are added in the output layer as demonstrated, e.g., on slide 15:



• [-] An MLP with a single hidden layer may require an infinite number of nodes.

- Answer: it is possible, but a general construction is more complicated as you might think.
- In this way any arbitrary contour can be approximated with arbitrary precision, with an MLP with a single hidden layer by filling the contours with cylinders. The perceptrons of the hidden layer are added in the output layer as demonstrated, e.g., on slide 15:



- [-] An MLP with a single hidden layer may require an infinite number of nodes.
- [+] Also here depth can allow a significant reduction of the required nodes<sup>[2]</sup>. How many hidden nodes for the example on the left?

- Answer: it is possible, but a general construction is more complicated as you might think.
- In this way any arbitrary contour can be approximated with arbitrary precision, with an MLP with a single hidden layer by filling the contours with cylinders. The perceptrons of the hidden layer are added in the output layer as demonstrated, e.g., on slide 15:



- [-] An MLP with a single hidden layer may require an infinite number of nodes.
- [+] Also here depth can allow a significant reduction of the required nodes<sup>[2]</sup>. How many hidden nodes for the example on the left? – 12: 10 in the first + 2 in the second hidden layer.

# Summary

PRO: It is possible to approximate any arbitrary Boolean function (slides 9–14), boundary (slides 15–19), or real-valued function (equivalent to slides 15–19) to arbitrary precision with an <u>MLP with only one single hidden layer</u>.

# Summary

- PRO: It is possible to approximate any arbitrary Boolean function (slides 9–14), boundary (slides 15–19), or real-valued function (equivalent to slides 15–19) to arbitrary precision with an <u>MLP with only one single hidden layer</u>.
- CON: A solution with only one single hidden layer in general may require an infinite amount of nodes.

# Summary

- PRO: It is possible to approximate any arbitrary Boolean function (slides 9–14), boundary (slides 15–19), or real-valued function (equivalent to slides 15–19) to arbitrary precision with an <u>MLP with only one single hidden layer</u>.
- CON: A solution with only one single hidden layer in general may require an infinite amount of nodes.
- PRO: The number of required nodes can be significantly reduced when exploiting depth for an MLP → depth matters.
# Summary

- PRO: It is possible to approximate any arbitrary Boolean function (slides 9–14), boundary (slides 15–19), or real-valued function (equivalent to slides 15–19) to arbitrary precision with an <u>MLP with only one single hidden layer</u>.
- CON: A solution with only one single hidden layer in general may require an infinite amount of nodes.
- PRO: The number of required nodes can be significantly reduced when exploiting depth for an MLP → depth matters.



# Summary

- PRO: It is possible to approximate any arbitrary Boolean function (slides 9–14), boundary (slides 15–19), or real-valued function (equivalent to slides 15–19) to arbitrary precision with an <u>MLP with only one single hidden layer</u>.
- CON: A solution with only one single hidden layer in general may require an infinite amount of nodes.
- PRO: The number of required nodes can be significantly reduced when exploiting depth for an MLP → depth matters.



# Summary

- PRO: It is possible to approximate any arbitrary Boolean function (slides 9–14), boundary (slides 15–19), or real-valued function (equivalent to slides 15–19) to arbitrary precision with an <u>MLP with only one single hidden layer</u>.
- CON: A solution with only one single hidden layer in general may require an infinite amount of nodes.
- PRO: The number of required nodes can be significantly reduced when exploiting depth for an MLP → depth matters.



• Deep NNs are **more** *expressive* than shallow NNs.

Priv.-Doz. Dr. Roger Wolf http://ekpwww.physik.uni-karlsruhe.de/~rwolf/

• Check out the following example:



**NB**: Ignore the bounding box and assume the chess-board structure to extend to infinity beyond the box.

• How many hidden nodes do we need to represent this contour with an NN with a <u>single</u> <u>hidden layer</u>?

• Check out the following example:



**NB**: Ignore the bounding box and assume the chess-board structure to extend to infinity beyond the box.

• How many hidden nodes do we need to represent this contour with an NN with a <u>single</u> <u>hidden layer</u>? – infinitely many!



**NB**: Ignore the bounding box and assume the chess-board structure to extend to infinity beyond the box.

- How many hidden nodes do we need to represent this contour with an NN with a <u>single</u> <u>hidden layer</u>? infinitely many!
- How many hidden nodes do we need to represent this contour with an NN with <u>two hidden</u> <u>layers</u>?



**NB**: Ignore the bounding box and assume the chess-board structure to extend to infinity beyond the box.

- How many hidden nodes do we need to represent this contour with an NN with a <u>single</u> <u>hidden layer</u>? infinitely many!
- How many hidden nodes do we need to represent this contour with an NN with <u>two hidden</u> <u>layers</u>? – 61! How do we get to this number?

• Check out the following example:





• Check out the following example:





• Check out the following example:





• Check out the following example:





• Check out the following example:





• Check out the following example:





• Check out the following example:





• Check out the following example:





• Check out the following example:





• Check out the following example:





• Check out the following example:





• Check out the following example:



• The contour consists of 12 lines, which can be represented by 12 nodes in the first hidden layer.



Priv.-Doz. Dr. Roger Wolf http://ekpwww.physik.uni-karlsruhe.de/~rwolf/

• Check out the following example:



• Check out the following example:



• The contour consists of 12 lines, which

Priv.-Doz. Dr. Roger Wolf http://ekpwww.physik.uni-karlsruhe.de/~rwolf/

• Check out the following example:



- Check out the following example:
- can be represented by 12 nodes in the  $x_2$ first hidden layer. Value 1 Value 0 • In the second hidden layer  $7 \times 7 = 49$ boxes are combined from the output of the first layer.  $x_1$ 49 nodes12 nodes $\sum : 12 + 49 = 61$  $x_2$

- Check out the following example:
- can be represented by 12 nodes in the  $x_2$ first hidden layer. Value 1 Value 0 • In the second hidden layer  $7 \times 7 = 49$ boxes are combined from the output of the first layer.  $x_1$ 49 nodes12 nodes**NB**: Removing only one single node from one of the layers means that again an  $x_2$ infinite number of nodes is required and/or the contour can only be approximated.

# Whish meets reality



- Imagine you don't have 12 but only 6 nodes at hand in the first layer, which you arrange as shown on the left.
- Also this arrangement requires an infinite number of nodes. Why?

# Whish meets reality



- Imagine you don't have 12 but only 6 nodes at hand in the first layer, which you arrange as shown on the left.
- Also this arrangement requires an infinite number of nodes. Why? – After passing the first layer you don't know any more where, e.g. inside an orange box, a sample is located.

# Whish meets reality



- Imagine you don't have 12 but only 6 nodes at hand in the first layer, which you arrange as shown on the left.
- Also this arrangement requires an infinite number of nodes. Why? – After passing the first layer you don't know any more where, e.g. inside an orange box, a sample is located.
- Root cause is the use of a step function as activation function in our example.

• Check out the following example:



- Imagine you don't have 12 but only 6 nodes at hand in the first layer, which you arrange as shown on the left.
- Also this arrangement requires an infinite number of nodes. Why? – After passing the first layer you don't know any more where, e.g. inside an orange box, a sample is located.
- Root cause is the use of a step function as activation function in our example.

Choosing an activation function for the nodes that preserves the information where, relative to the hyperplane boundary, a sample is located (slide 4) allows **processing of information in proceeding layers**. This works the better, the better the activation function supports the transmission of this information through consecutive layers.



## Machine learning



• Up to now we did not touch the actual ML part: the question how to determine the weights and thresholds of the MLPs to fulfill their tasks.

# <sup>30/40</sup> Truth vs. prediction

 Assume the MLP should represent the blue function shown on the right (→ truth).



## Truth vs. prediction

- Assume the MLP should represent the blue function shown on the right (→ truth).
- Random choice of the weights {ω<sub>ij</sub>} might result in the red curve, shown on the right (→ prediction).



## Truth vs. prediction

- Assume the MLP should represent the blue function shown on the right (→ truth).
- Random choice of the weights {ω<sub>ij</sub>} might result in the red curve, shown on the right (→ prediction).
- **Task:** Adapt the weights such that the red curve approaches the blue one as closely as possible.



## Truth vs. prediction

- Assume the MLP should represent the blue function shown on the right (→ truth).
- Random choice of the weights {ω<sub>ij</sub>} might result in the red curve, shown on the right (→ prediction).
- **Task:** Adapt the weights such that the red curve approaches the blue one as closely as possible.
- To solve this task mathematically we will quantify the difference between the two curves with the socalled **loss or cost function**.



## Sample and training

• In general we don't know the truth. We have to infer it from a sample hoping that the sample is *representative* of the ground truth (→ learning by example).

## Sample and training

- In general we don't know the truth. We have to infer it from a sample hoping that the sample is *representative* of the ground truth (→ learning by example).
- Learning by example we call training.
# Sample and training

- In general we don't know the truth. We have to infer it from a sample hoping that the sample is *representative* of the ground truth (→ learning by example).
- Learning by example we call training.
- To be representative the sample should catch all relevant characteristics of the truth. Individual properties of the sample (a.k.a. fluctuations) should not influence the training (→ generalization property).

# Sample and training

- In general we don't know the truth. We have to infer it from a sample hoping that the sample is *representative* of the ground truth (→ learning by example).
- Learning by example we call **training**.
- To be representative the sample should catch all relevant characteristics of the truth. Individual properties of the sample (a.k.a. fluctuations) should not influence the training (→ generalization property).
- If an MLP has so many trainable parameters that it can pick up on fluctuations of the sample it runs into the issue of overfitting (→ overtraining). If the MLP has too few parameters it might not be *expressive* enough to do the job (→ underfitting).

• **Historic example**: Training of a single Boolean perceptron to separate two classes with the help of labeled examples (here represented by points with different color)<sup>(1)</sup>:



Task: Determine the weights of the perceptron such that the red points (with values 1) and the blue points (with values 0) are separated.

<sup>(1)</sup> This is the historic example by Frank Rosenblatt, from the 1960ies.

# <sup>32/40</sup> The perceptron learning rule

• **Historic example**: Training of a single Boolean perceptron to separate two classes with the help of labeled examples (here represented by points with different color)<sup>(1)</sup>:



- Task: Determine the weights of the perceptron such that the red points (with values 1) and the blue points (with values 0) are separated.
- **NB**: The solution is ambiguous. The root cause of this is that the sample does not cover the complete space over  $x_1$  and  $x_2$ .

<sup>(1)</sup> This is the historic example by Frank Rosenblatt, from the 1960ies.

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



$$\boldsymbol{\omega} \cdot \mathbf{x} < 0 \qquad --- \quad \boldsymbol{\omega} \cdot \mathbf{x} = 0 \qquad \boldsymbol{\omega} \cdot \mathbf{x} > 0$$

- Initialize weights randomly.
- Only update for examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} 
ightarrow \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} 
ight.$$

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



- Initialize weights randomly. Step 0
  - Only update for examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \to \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} 
ight.$$

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



- Initialize weights randomly.
- Only update for examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \to \boldsymbol{\omega}^{(k+1)} = \begin{cases} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & \text{if red} \\ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & \text{if blue} \end{cases}$$
Step 1

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



- Initialize weights randomly.
- Only update for examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \to \boldsymbol{\omega}^{(k+1)} = \begin{cases} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & \text{if red} \\ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & \text{if blue} \end{cases}$$
Step 1

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



- Initialize weights randomly.
- Only update for examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \rightarrow \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} 
ight.$$

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



**NB**: In the backup you can find the same algorithm played through w/ the blue points.

### Algorithm:

- Initialize weights randomly.
- Only update for examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \to \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} 
ight.$$

• **NB**: Rosenblatt could show that a single logic perceptron for linearly separable tasks always converges to the correct solution after a *finite number* of steps.

### The perceptron learning rule – limitations

• The attempt to apply the same learning rule to more complex problems fails. Q: Why?



- Initialize weights randomly.
- Only update for examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} 
ightarrow \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} 
ight.$$

## The perceptron learning rule – limitations

- The attempt to apply the same learning rule to more complex problems fails. **Q**: Why?
- A: For the given example the samples are **not linearly separable**, i.e. one cannot draw a line to separate the blue from the red points!



- Initialize weights randomly.
- Only update for examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \to \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} 
ight.$$

### The perceptron learning rule – limitations

- The attempt to apply the same learning rule to more complex problems fails. **Q**: Why?
- A: For the given example the samples are **not linearly separable**, i.e. one cannot draw a line to separate the blue from the red points!



### Algorithm:

- Initialize weights randomly.
- Only update for examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \to \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} 
ight.$$

• **NB**: The problem can still be solved. One could re-label the values 0 and 1 for each hyperplane, but the complexity of the problem would still grow exponentially.

# The role of the activation function – revisited

- The issue of the perceptron learning rule originates from the use of the Heavyside step function as activation function:
- Small changes of the weights of the perceptron have **no influence** on its output.



# The role of the activation function – revisited

- The issue of the perceptron learning rule originates from the use of the Heavyside step function as activation function:
- Small changes of the weights of the perceptron have **no influence** on its output.



- **Solution**: Introduce a continuous activation function.
- This turns the NN output function differentiable in <u>any</u> variable:

$$y^{(k)}(\boldsymbol{\omega}, \mathbf{x}) = \theta(z(\boldsymbol{\omega}, \mathbf{x})) \quad z^{(k)} = \sum_{i} \omega_{i} y_{i}^{(k-1)}$$
$$\frac{\mathrm{d}y^{(k)}}{\mathrm{d}\omega_{i}} = \frac{\mathrm{d}y^{(k)}}{\mathrm{d}z^{(k)}} \frac{\mathrm{d}z^{(k)}}{\mathrm{d}\omega_{i}} = \theta'(z^{(k)}) y_{i}^{(k-1)}$$
$$\frac{\mathrm{d}y^{(k)}}{\mathrm{d}x_{i}} = \frac{\mathrm{d}y^{(k)}}{\mathrm{d}z^{(k)}} \frac{\mathrm{d}z^{(k)}}{\mathrm{d}x_{i}} = \theta'(z^{(k)}) \omega_{i}$$

• This allows checking what changes in  $y^{(k)}$  one obtains from small changes in  $\omega$ .

# The continuous activation function

- The issue of the perceptron learning rule originates from the use of the Heavyside step function as activation function:
- **Example** sigmoid function:



Wrongly predicted samples





• Gradual changes of the weights lead to gradual changes of the NN output.

• We take the sigmoid function as an example of addressing not linearly separable tasks, here demonstrated with a 1d example:





# <sup>39/40</sup> Not linearly separable tasks

• We take the sigmoid function as an example of addressing not linearly separable tasks, here demonstrated with a 1d example:





# <sup>39/40</sup> Not linearly separable tasks

• We take the sigmoid function as an example of addressing not linearly separable tasks, here demonstrated with a 1d example:





# <sup>39/40</sup> Not linearly separable tasks

• We take the sigmoid function as an example of addressing not linearly separable tasks, here demonstrated with a 1d example:





• We take the sigmoid function as an example of addressing not linearly separable tasks, here demonstrated with a 1d example:





### <sup>40/40</sup> Summary

- Continuous activation functions allow transmission of information through hyperplanes.
- Gradual changes of the weights lead to gradual changes of the NN output. In this way the minimization task can in principal be solved analytically!
- **NB**: To solve the minimization task you also need the differentiability of the loss function.

- For practical reasons neither the loss, nor the NN output function need to have a total derivative. It is sufficient if a gradual change in feature space leads to a gradual change in the NN output and loss function.
- An example of an activation function, which is not fully differentiable in all points of its input space is ReLU (slide 4).

### Backup

**Solution**: Hyperplane in Hessian canonical form  $\sum \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane • (=on the boundary).



### Algorithm:

Initialize weights randomly.



- Only update only examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \to \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} 
ight.$$

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



- Initialize weights randomly.
- Only update only examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \rightarrow \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \\ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} \middle| egin{array}{cc} \mathsf{Step 1} \end{array} 
ight.$$

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



- Initialize weights randomly.
- Only update only examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \rightarrow \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \\ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} \middle| egin{array}{cc} \mathsf{Step 1} \end{array} 
ight.$$

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



- Initialize weights randomly.
- Only update only examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \to \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \\ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} \right.$$
 Step 2

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



- Initialize weights randomly.
- Only update only examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \to \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \\ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} \right.$$
 Step 2

• Solution: Hyperplane in Hessian canonical form  $\sum_{i} \omega_i x_i = 0$  i.e.  $\omega \perp \mathbf{x} \quad \forall \mathbf{x}$  in the plane (=on the boundary).



- Initialize weights randomly.
- Only update only examples w/ wrong predictions.
- For those, apply the following learning rule:

$$\boldsymbol{\omega}^{(k)} \rightarrow \boldsymbol{\omega}^{(k+1)} = \left\{ egin{array}{cc} \boldsymbol{\omega}^{(k)} + \mathbf{x}^{(k)} & ext{if red} \ \boldsymbol{\omega}^{(k)} - \mathbf{x}^{(k)} & ext{if blue} \end{array} 
ight.$$

# Notes on the lecture

- <sup>[1]</sup> In the previous lecture we have demonstrated that for the pentagon **only five** and not infitely many notes are necessary! This was the case, since for the pentagon the corners are connected by straight lines. In the slides following this annotation a general solution is discussed.
- <sup>[2]</sup> The difference to [1] is that here there is no alternative to filling the disconnected contours with dots.