



# WIE MAN EINE KI IN DIE IRRE FÜHRT

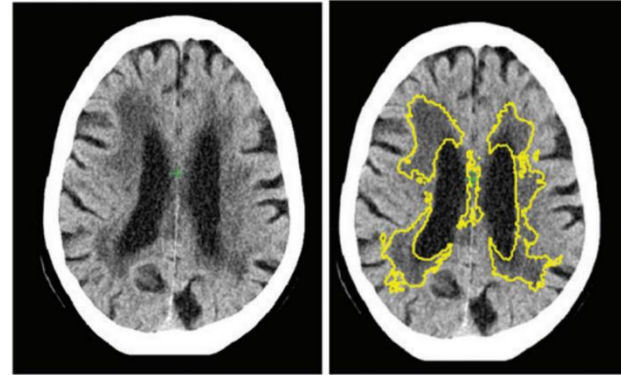
Anton Winschel

7. November 2019 – IT Kongress Neu-Ulm | Ulm

Autonomes Fahren



Medizinische Diagnosen



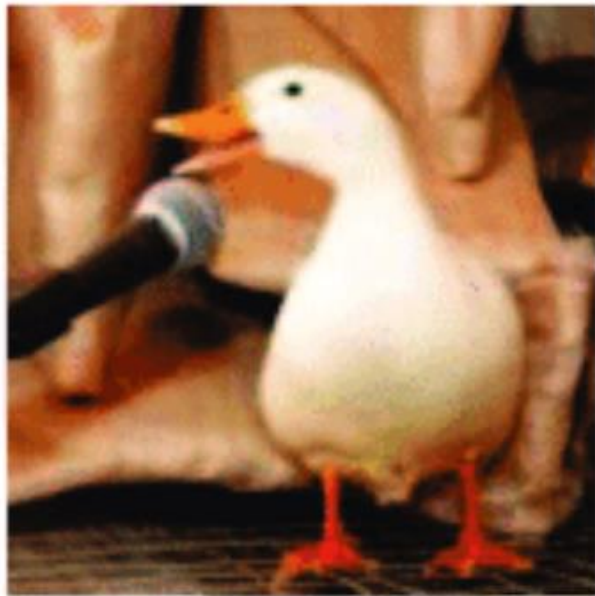
Industrie



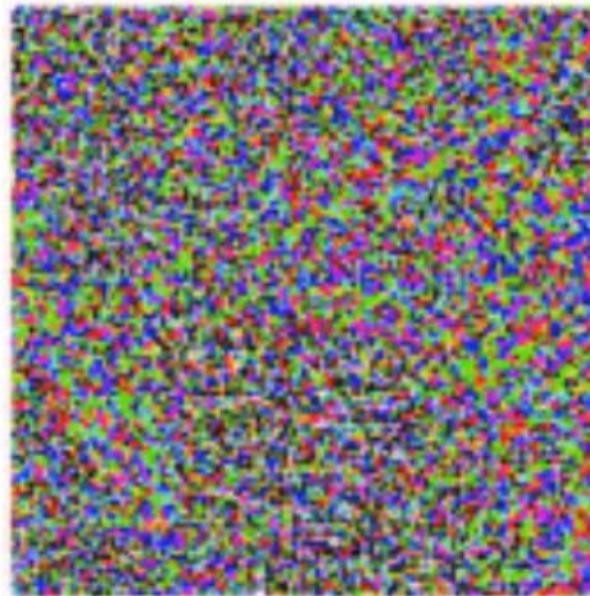
- › Wie kann man eine KI in die Irre führen?
- › Beispiel: Klassifikation
  - › Funktionsweise
  - › Irreführung
- › Wie können KI-Systeme robuster gestaltet werden?

Quelle: [medium.com/intro-to-artificial-intelligence](https://medium.com/intro-to-artificial-intelligence)  
[pubs.rsna.org/doi/pdf/10.1148/radiol.2018171567](https://pubs.rsna.org/doi/pdf/10.1148/radiol.2018171567)  
[www.cross-compass.com/article2](https://www.cross-compass.com/article2)

# IRREFÜHRUNG: OBJEKTERKENNUNG



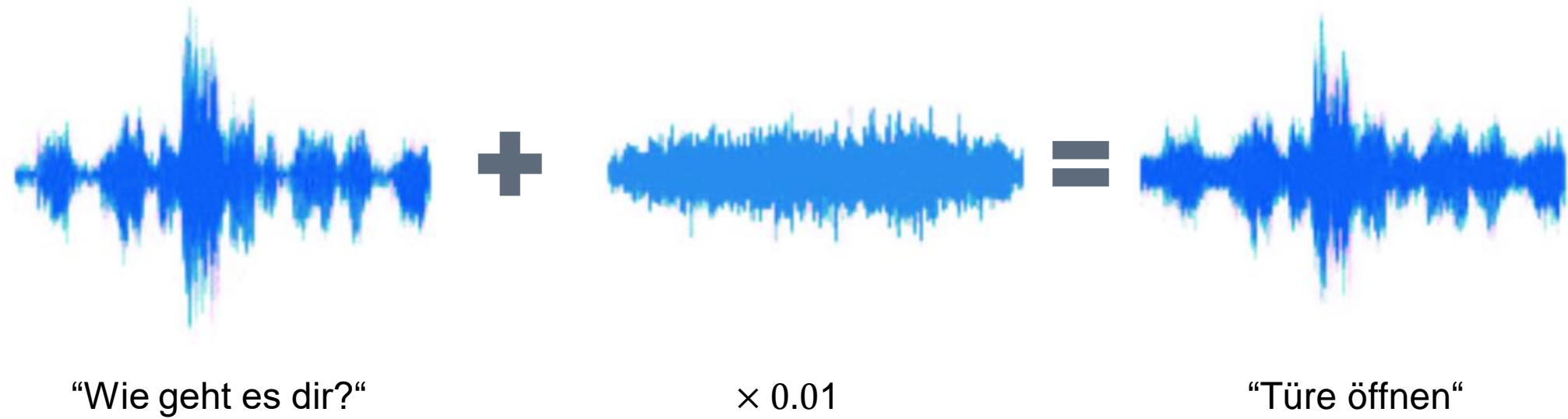
Ente



$\times 0.07$

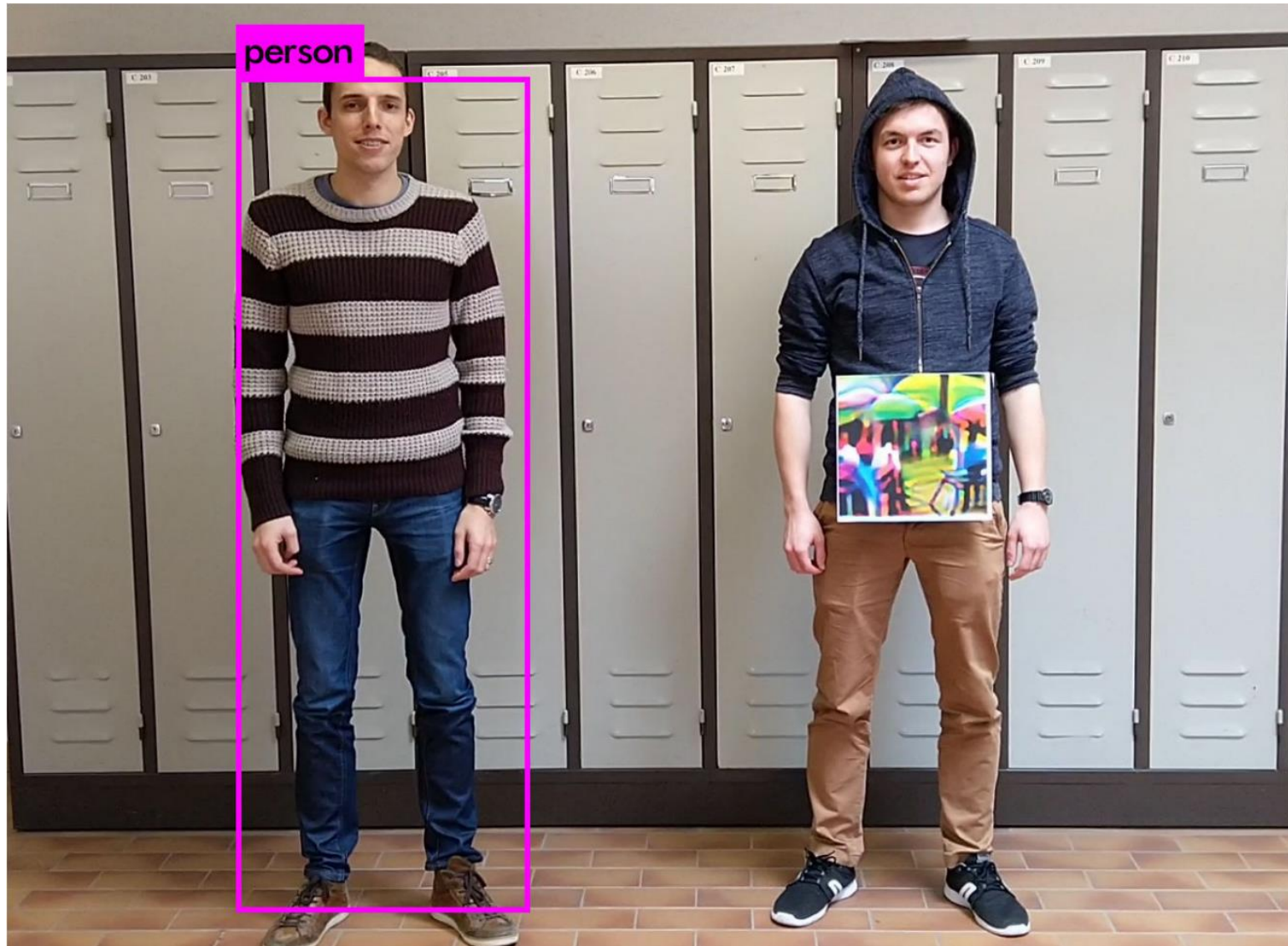


Pferd





# IRREFÜHRUNG: ERKENNUNG VON PERSONEN



“Unsichtbarkeits-T-Shirt”

Quelle: arXiv:1904.08653v1  
cloakwear.co



Limit 60



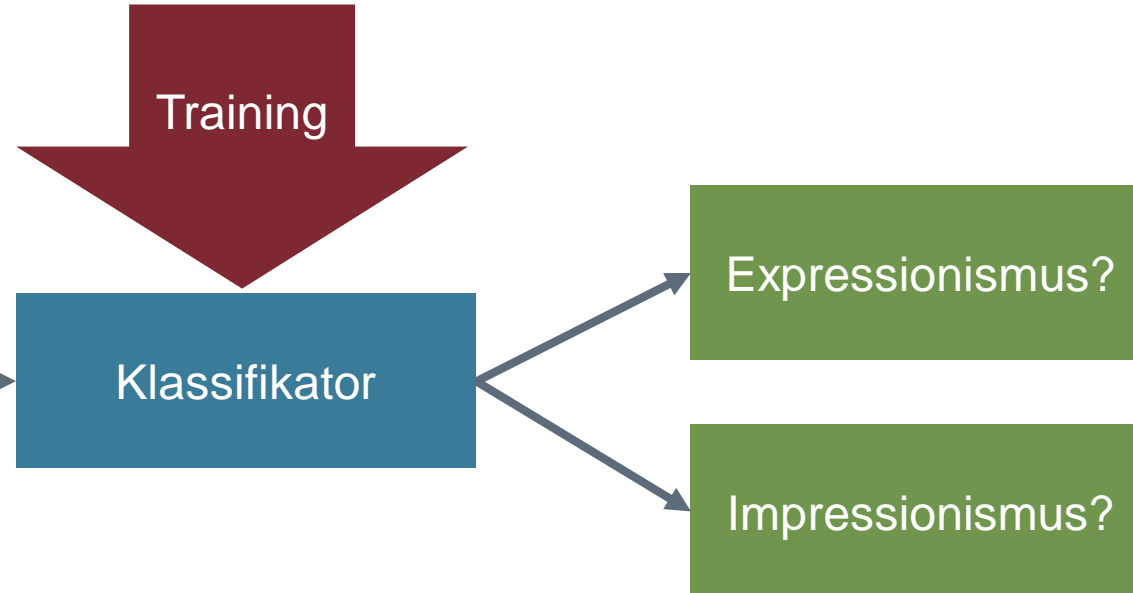
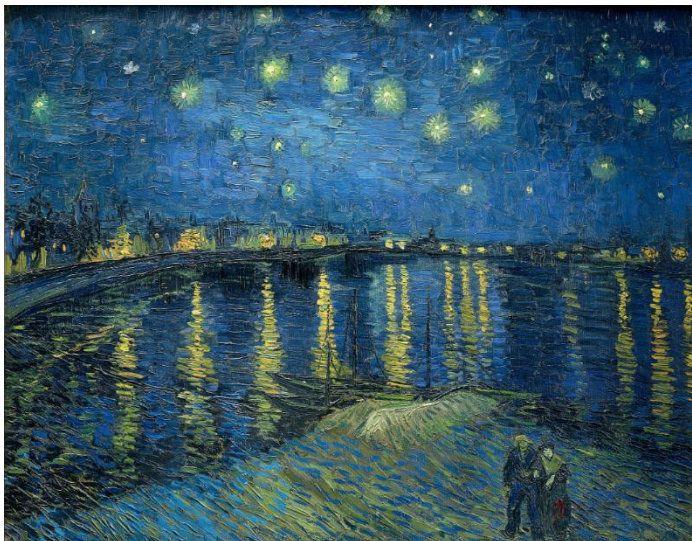
Limit 45

Quelle: arXiv:1707.08945v5

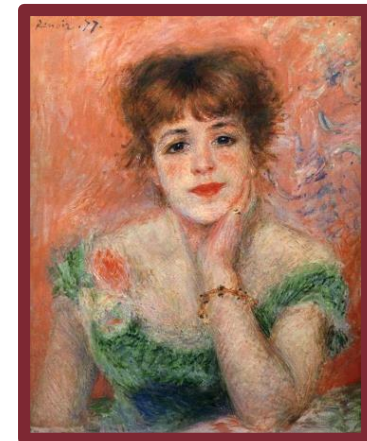
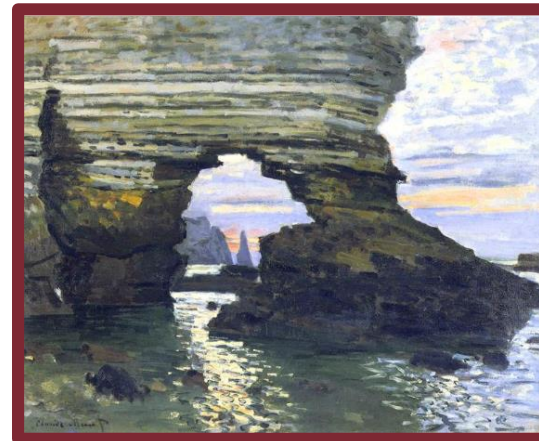
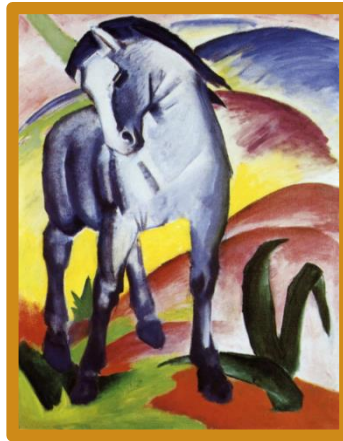
[sites.nlsde.buaa.edu.cn/~xlliu/aaai19.pdf](http://sites.nlsde.buaa.edu.cn/~xlliu/aaai19.pdf)



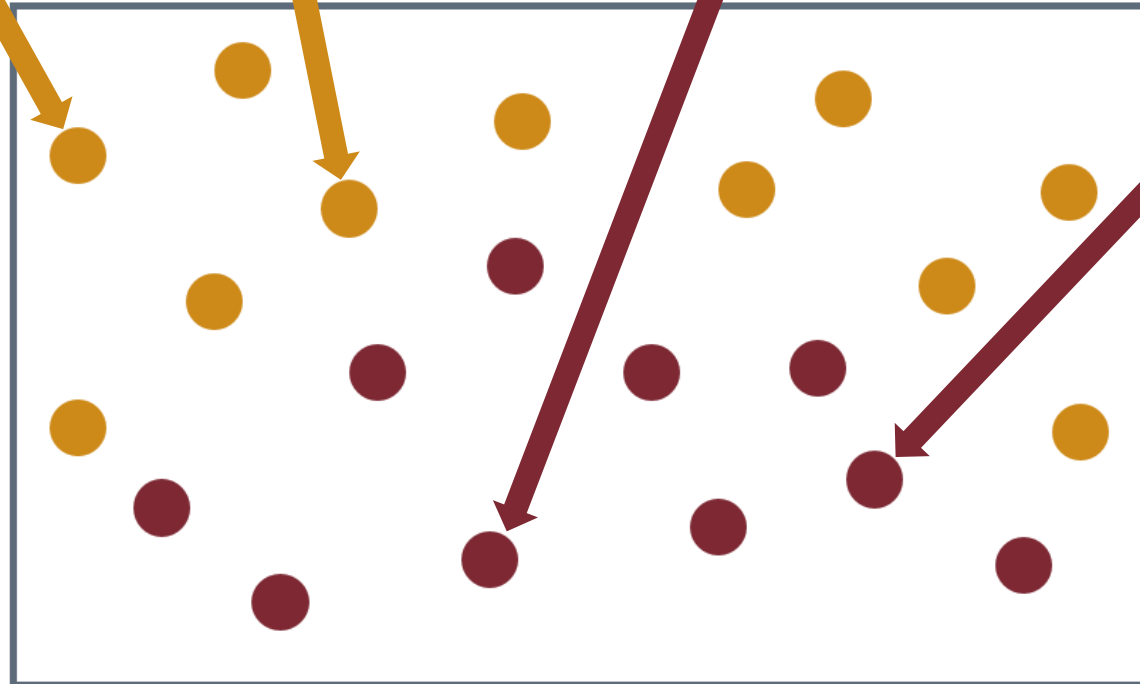
# KLASSIFIKATION – KUNSTRICHTUNG ERKENNEN



# KLASSIFIKATION – GEOMETRISCHE SICHT



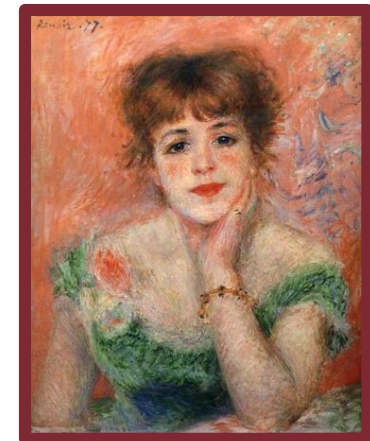
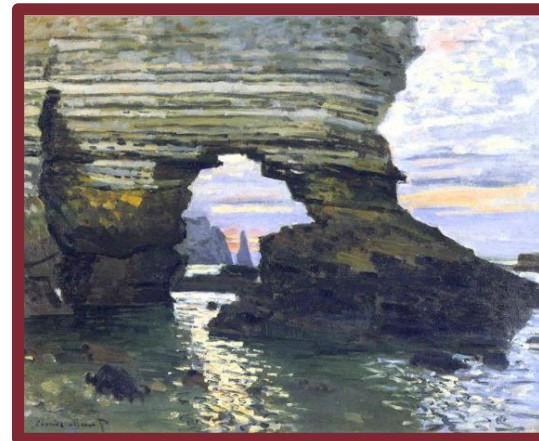
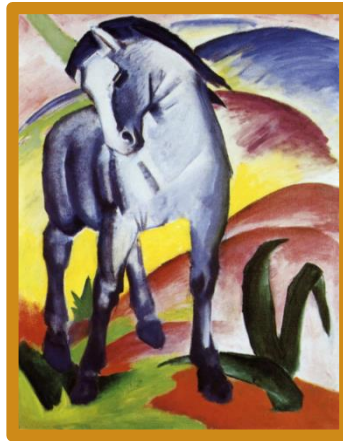
Jedes Bild wird als  
Punkt in einem  
hochdimensionalen  
Raum abgebildet



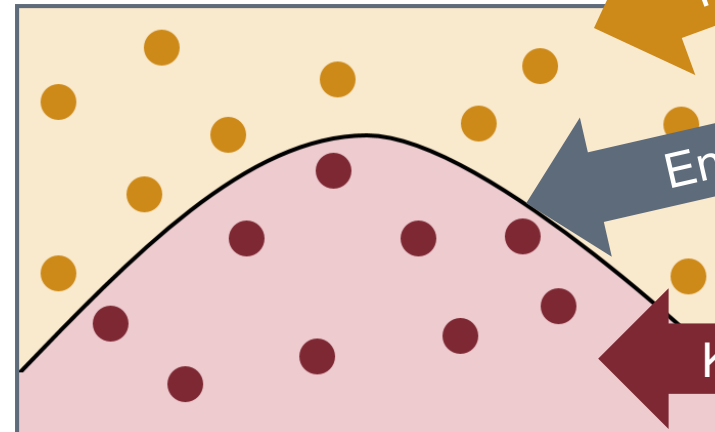
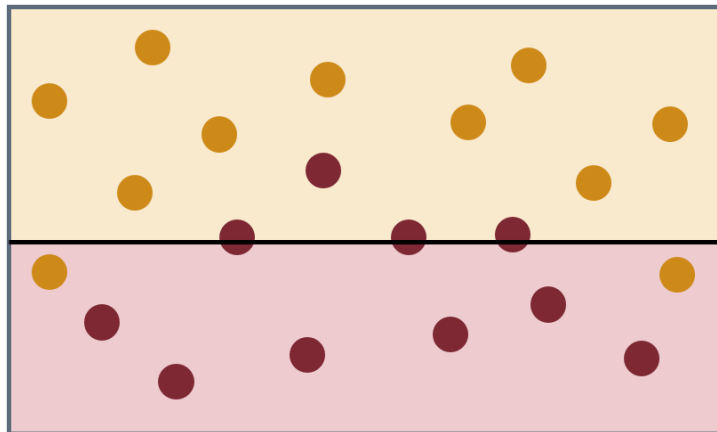
Z.B.  
 $500px \times 500px \times 3 (rgb)$   
 $\approx 700.000 \text{ dim}$



# KLASSIFIKATION – GEOMETRISCHE SICHT



Beim Training wird die Entscheidungsgrenze schrittweise angepasst

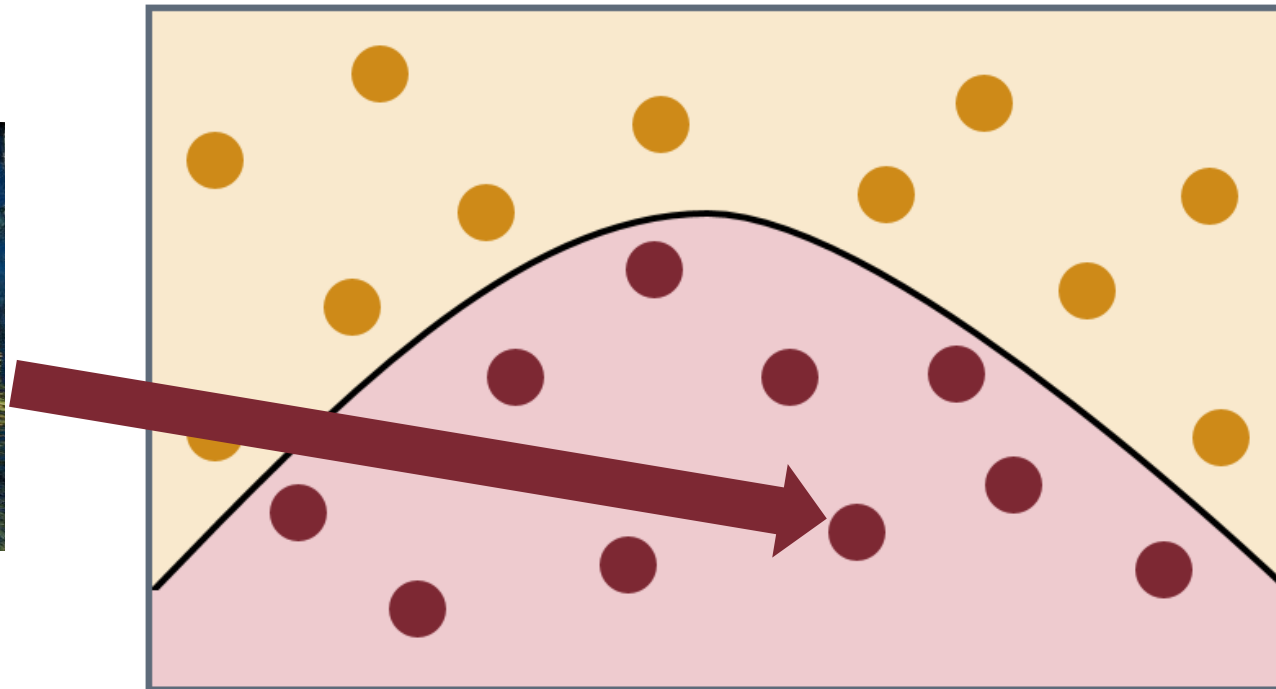
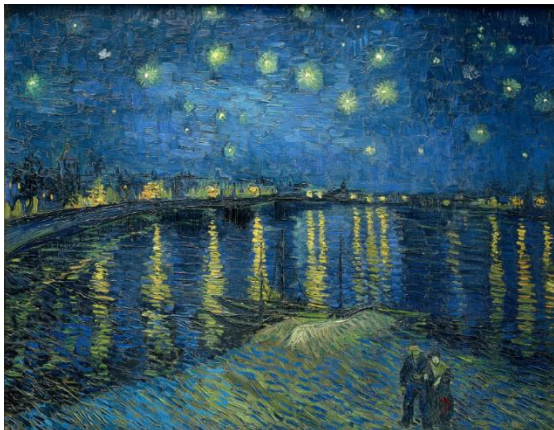
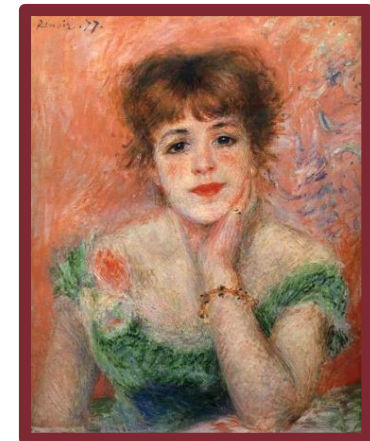
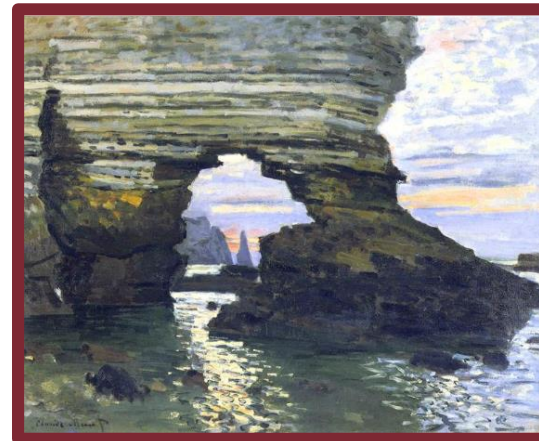
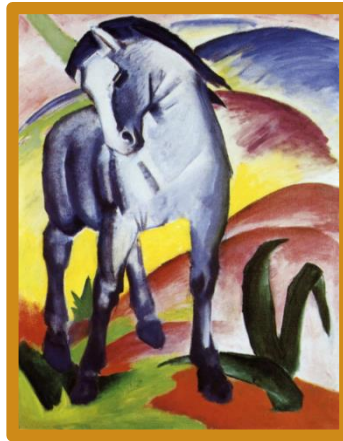


Klasse: Expressionismus

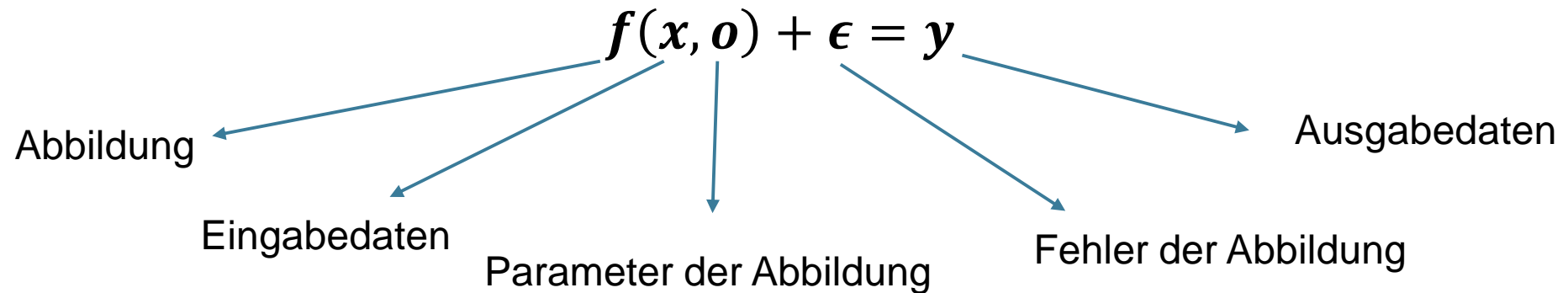
Entscheidungsgrenze

Klasse: Impressionismus

# KLASSIFIKATION – GEOMETRISCHE SICHT



Mit der  
Entscheidungsgrenze  
wird eine Klassifikation  
durchgeführt



$$f(\text{Expressionismus}, o) + \epsilon = \text{Expressionismus}$$

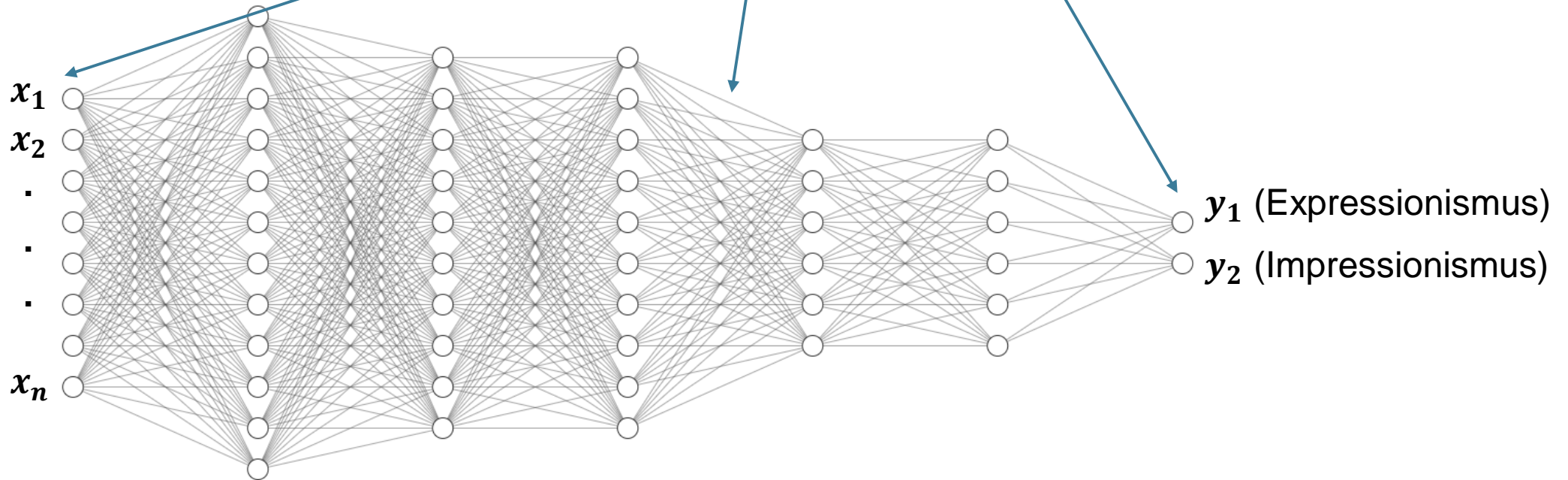
Ziel:

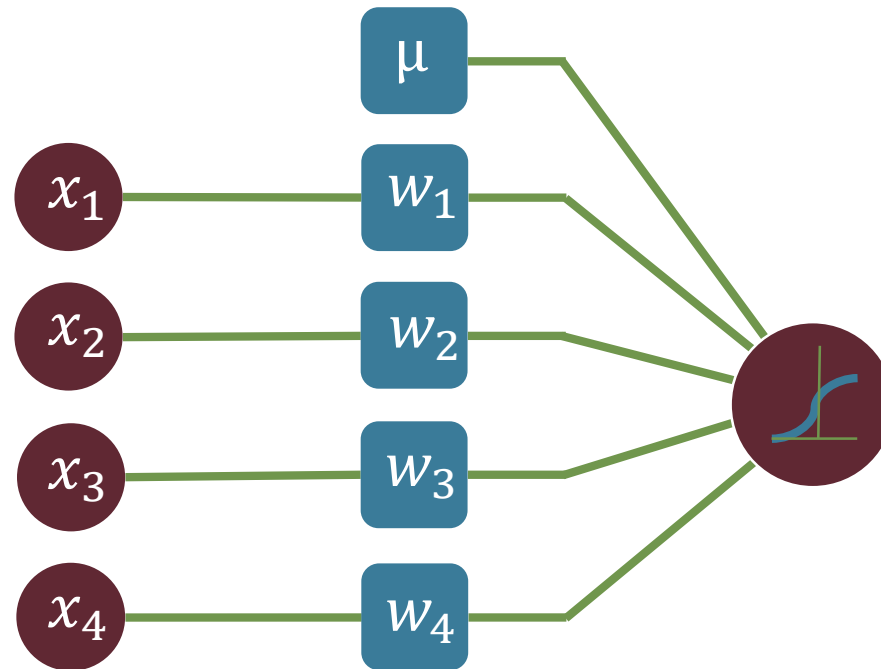
1. Finde  $o$ , sodass  $\epsilon$  minimal ist (Optimierungsproblem)
2.  $f$  soll auf unbekannte Daten generalisieren



Beispiel:  $f$  = Neuronales Netz

$$f(x, o) + \epsilon = y$$





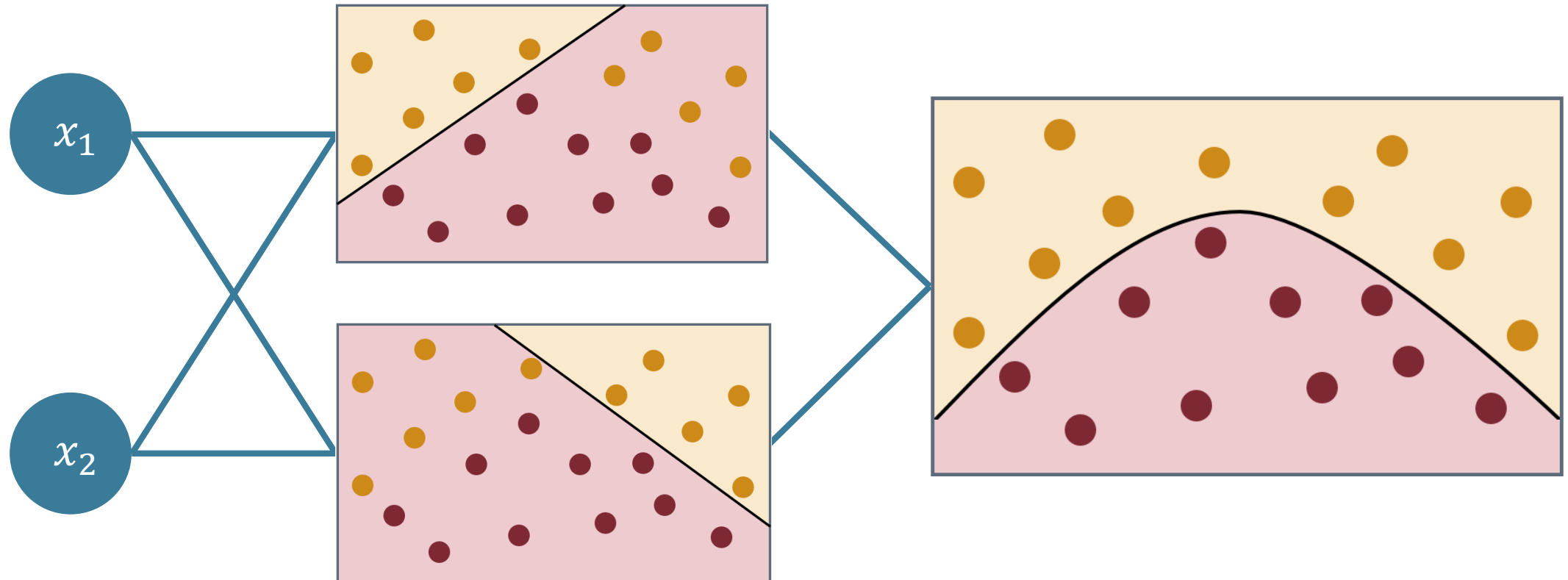
$$\sum_i x_i w_i - \mu = y$$

Entscheidungsgrenze

Jedes Neuron berechnet ein Skalarprodukt der Vektoren  $\mathbf{x}$  und  $\mathbf{w}$

# KLASSIFIKATION – GEOMETRISCHE SICHT

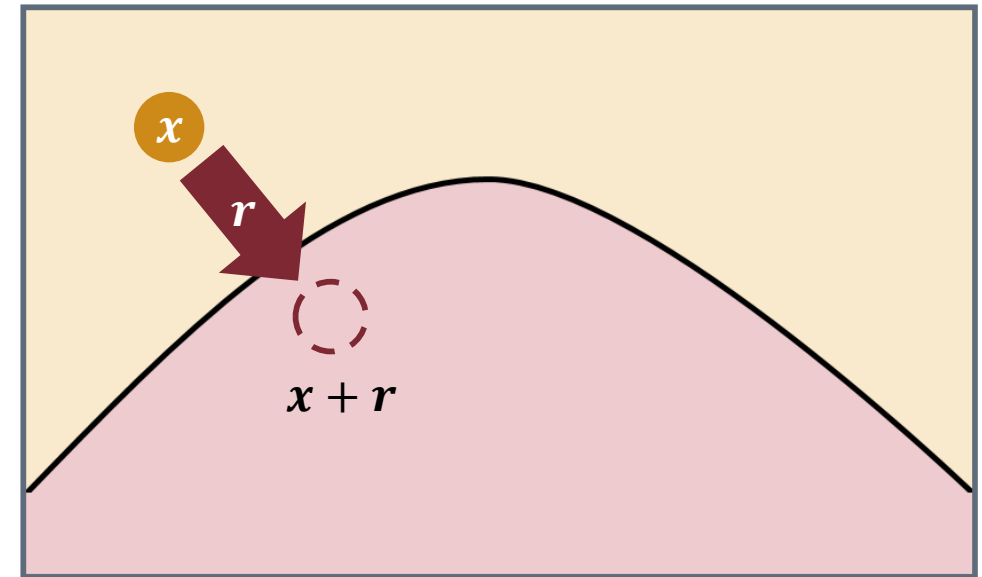
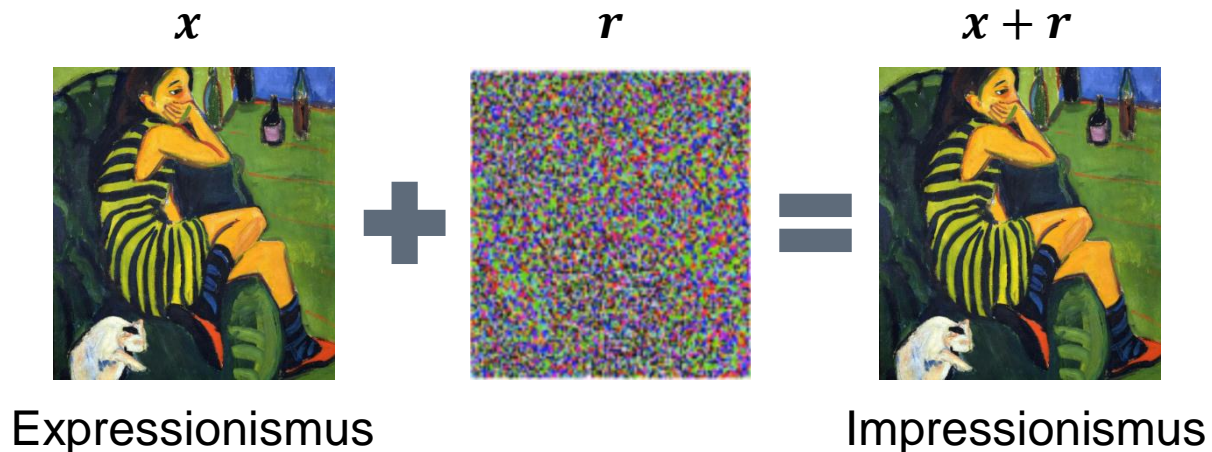
- › Jedes Neuron teilt den Eingangsraum in 2 Teile
- › Jedes Neuron ist ein binärer linearer Klassifikator
- › Ein neuronales Netz mit nur einer Schicht kann jede Funktion beliebig genau approximieren

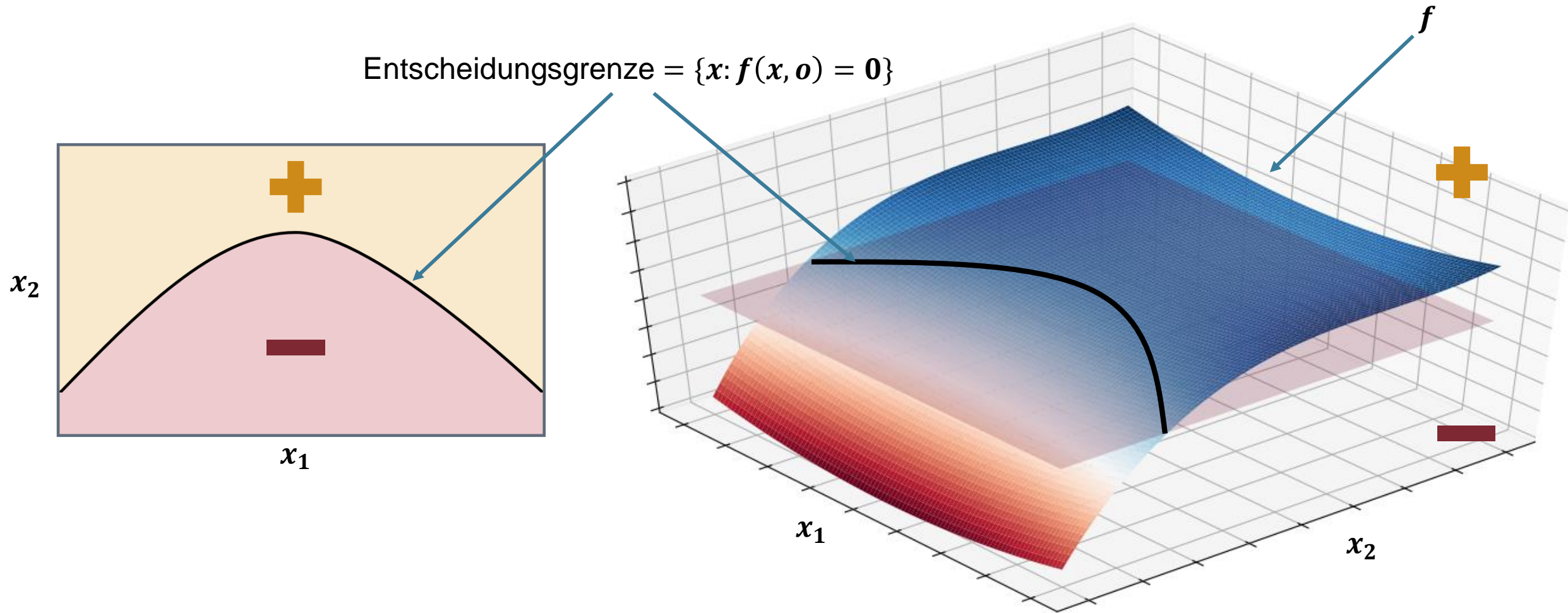


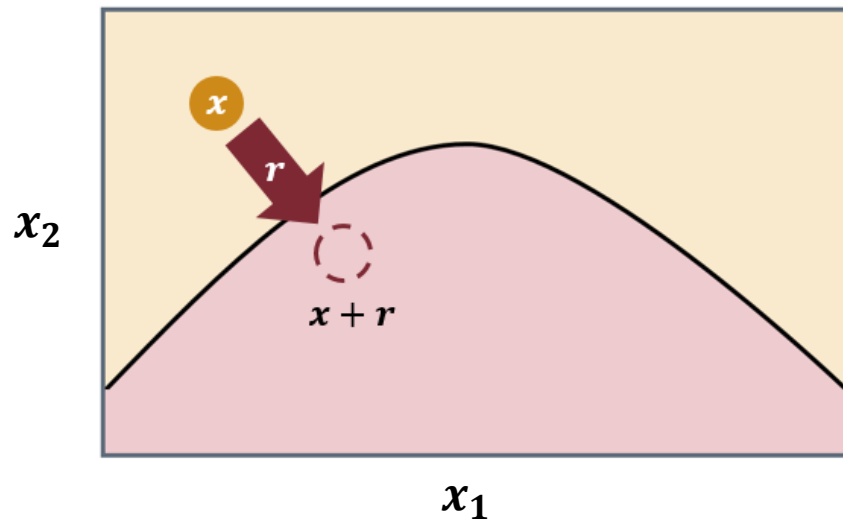


Verschiebung von  $x$  auf die andere Seite der Entscheidungsgrenze

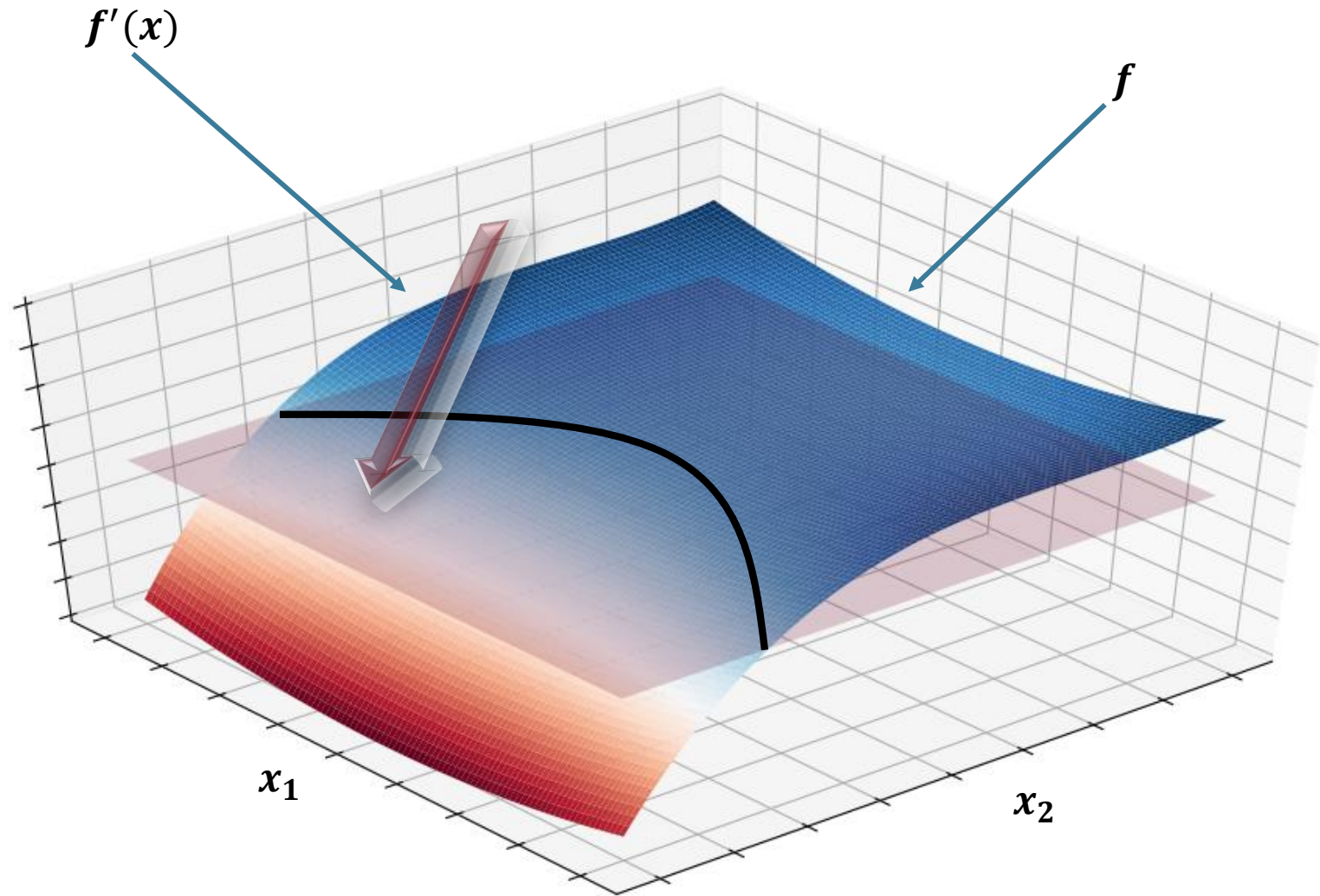
- › Klassifizierung ändert sich
- › Bei ungünstiger Entscheidungsgrenze ist kein visueller Unterschied erkennbar







$r$  kann durch die Ableitung von  $f$  berechnet werden





# IRREFÜHRUNG – BEISPIEL

- › Klassifikator für 1000 Klassen
- › Wir möchten dass ein Wal als Hai klassifiziert wird



Wal



Hai?

# IRREFÜHRUNG – BEISPIEL

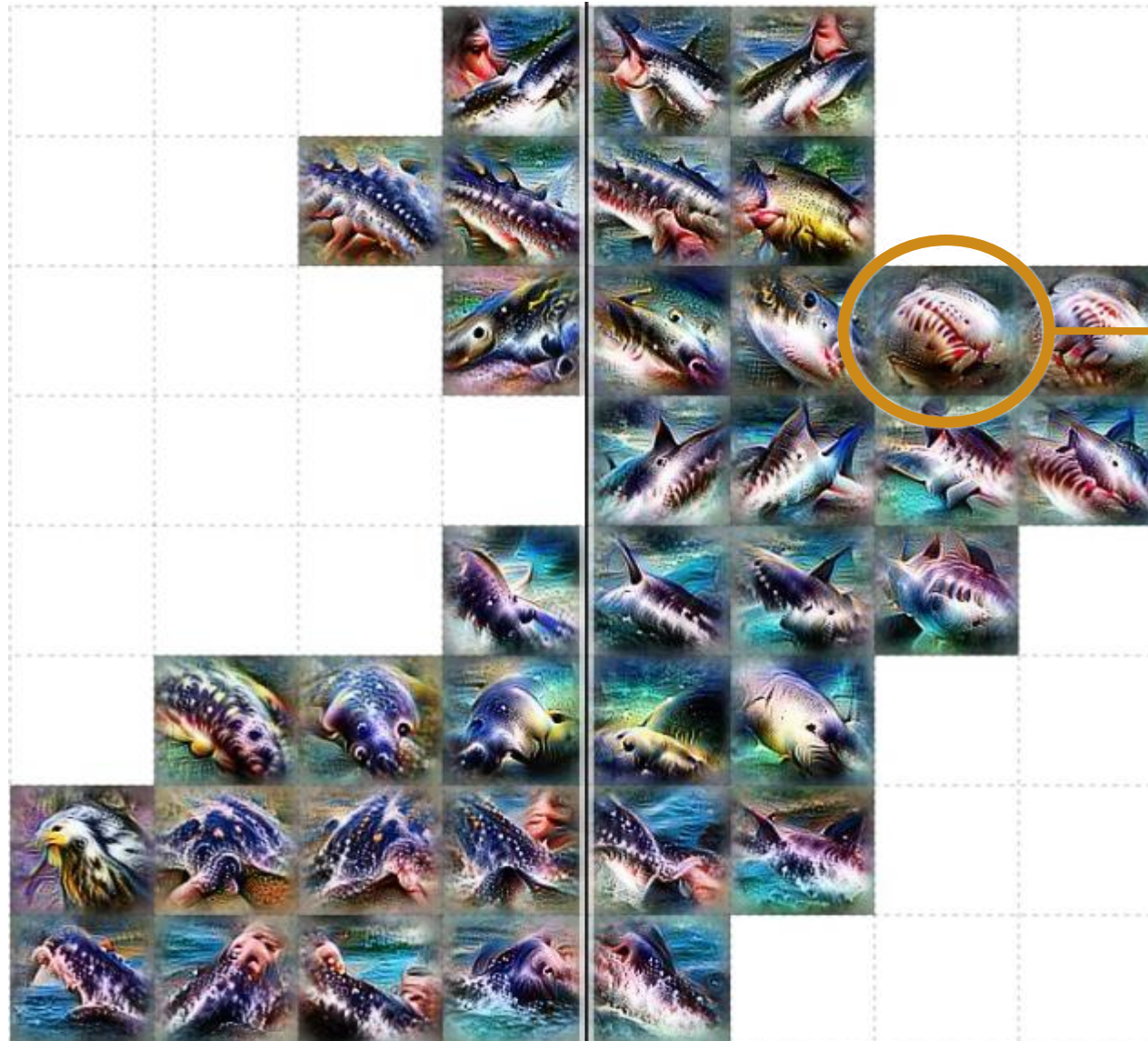
Entscheidungsgrenze

Wal

Hai

Baseball?

“Debugging“ von  
neuronalen Netzen



# IRREFÜHRUNG – BEISPIEL

$x$



Wal

$r$



Baseball

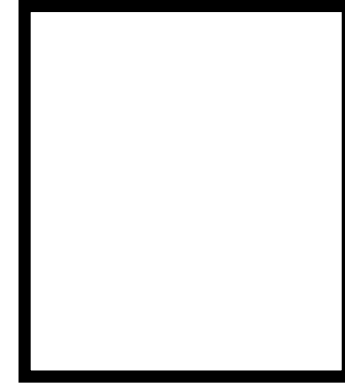
$x + r$



Hai



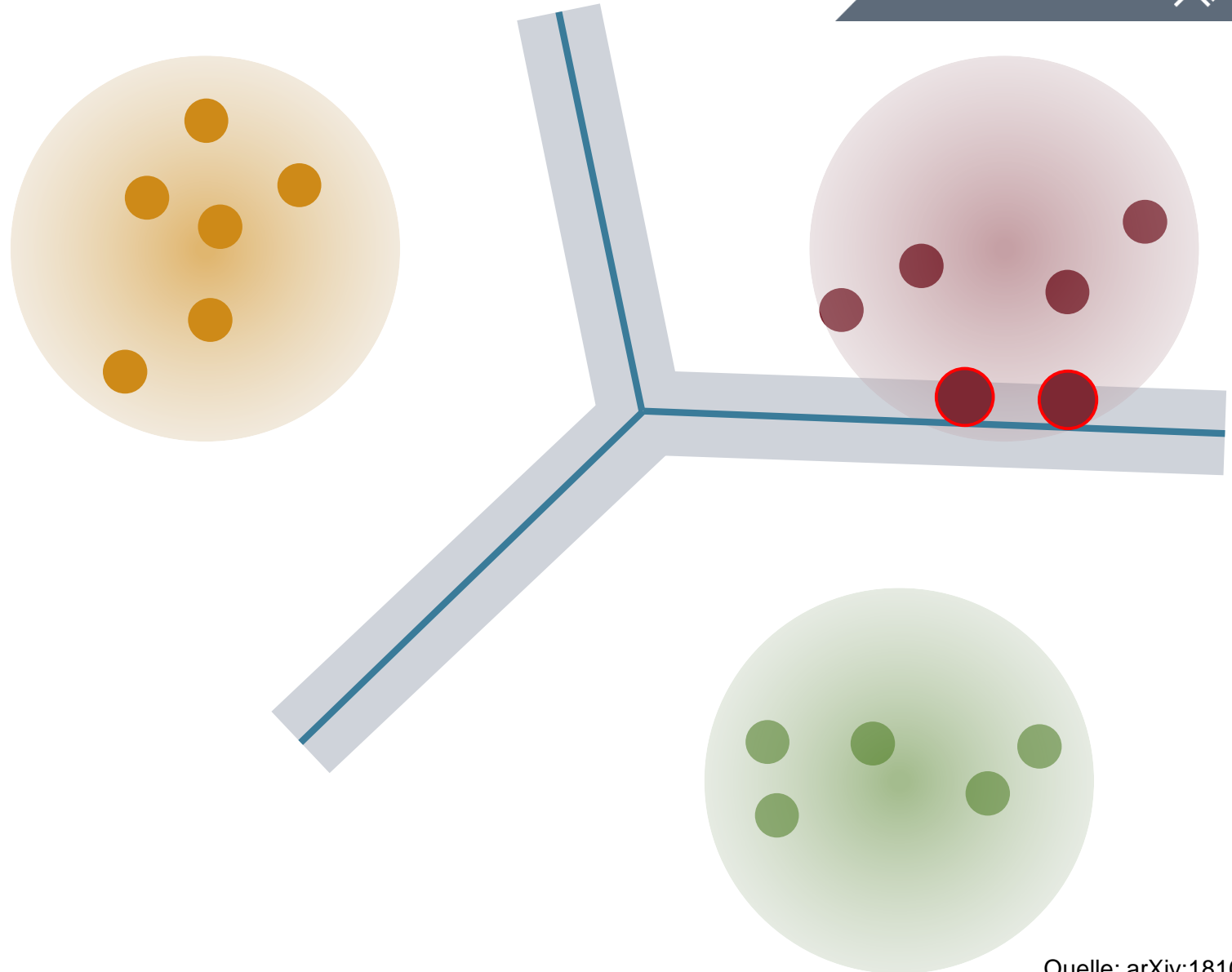
## Trainingsdatensatz



$$\begin{array}{c} x \\ \text{[Image of woman in striped dress]} \end{array} + \begin{array}{c} r \\ \text{[Image of random noise]} \end{array} = \begin{array}{c} x + r \\ \text{[Image of woman in striped dress with noise]} \end{array}$$

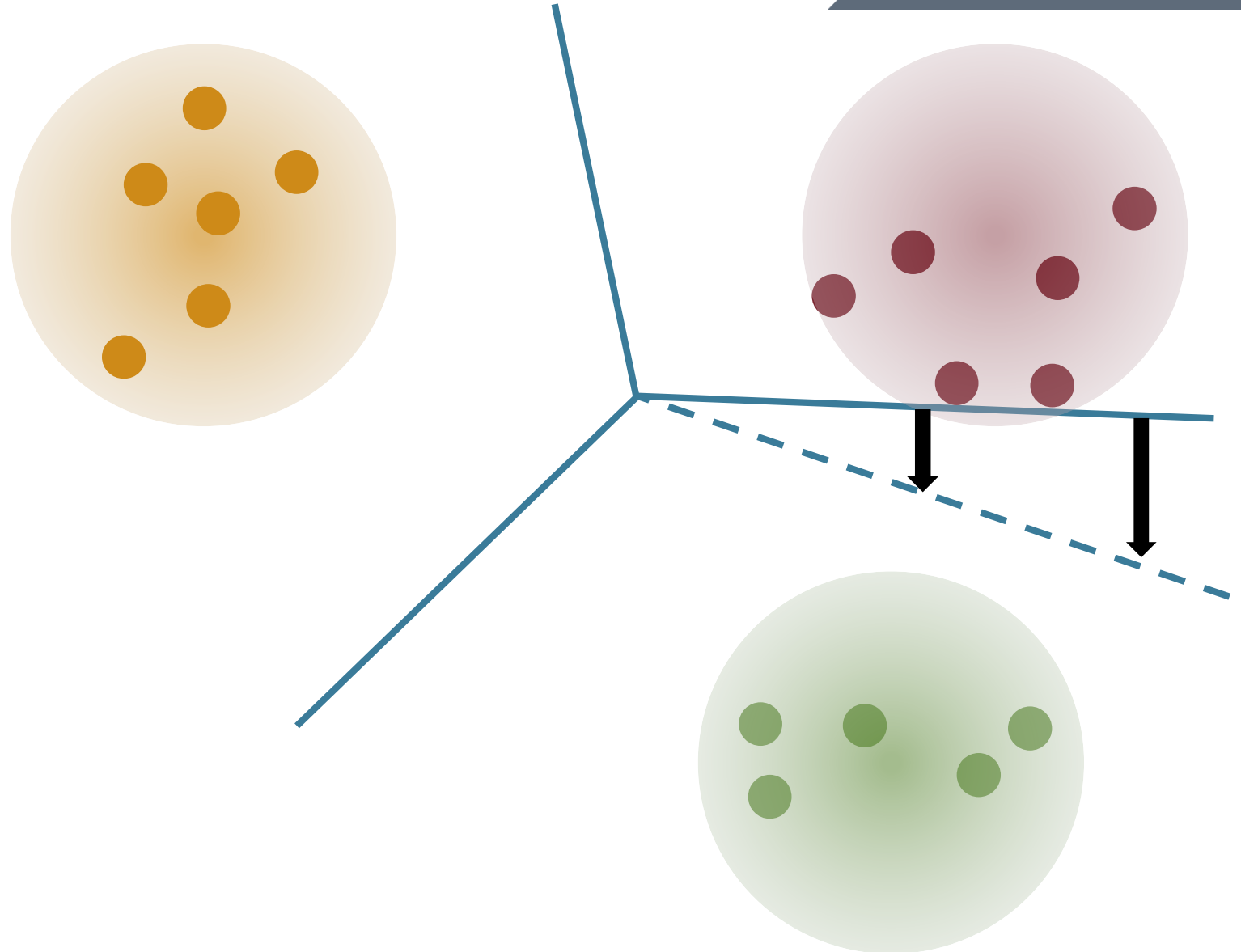


Irreführende Daten zum Training hinzunehmen  
Aber: Kein Schutz vor anderen Arten von Angriffen



Sonderbehandlung für  
Daten in der Nähe von  
Entscheidungsgrenzen

- › Entscheidungsgrenzen verschieben
- › Maximalen Abstand zu den Klassen gewährleisten
- › Aufgrund hoher Dimension und Nichtlinearität schwierig
- › Erste Lösungsansätze aus der Forschung: [arXiv:1810.12715v4](https://arxiv.org/abs/1810.12715v4)



- › KI wird bereits in vielen Szenarien sinnvoll eingesetzt
- › KI birgt aber auch Schwächen und Risiken (Robustheit, Angriffssicherheit, ...)
- › Bei einer nicht-robusten KI reicht zufälliges Sensorrauschen für eine Irreführung aus
- › Ein “normales” Training erzeugt kein robustes System
- › Für den robusten Einsatz in sicherheitskritischen Szenarien sind zusätzliche Maßnahmen nötig



## VIELEN DANK

**XITASO**  
in Augsburg

Austraße 35  
86153 Augsburg

Tel. +49 (0)821 885 882 0  
E-Mail [info@xitaso.com](mailto:info@xitaso.com)  
Web [www.xitaso.com](http://www.xitaso.com)

**XITASO**  
in Magdeburg

Werner-Heisenberg-Straße 1  
39106 Magdeburg

Tel. +49 (0)391 / 792 930 00  
E-Mail [info@xitaso.com](mailto:info@xitaso.com)  
Web [www.xitaso.com](http://www.xitaso.com)