

# Grundlagen der Künstlichen Intelligenz

## Wintersemester 25/26

### Vorlesung 3

## Maschinelles Lernen I: Lernende Algorithmen, Klassifikation & Logistische Regression

Prof. Dr. Pascal Friederich  
T.T.-Prof. Dr. Peer Nowack

# KI-Landkarte

## Künstliche Intelligenz

### Modellierung und Schlussfolgerung

Variablen VL12 Inferenz

Logik Wissensrepräsentation

VL11

Zustände Suche VL13 MDPs

Reflex

### Anwendungen

Robotik VL14

Computer Vision VL9

Natürliche Sprache VL10

### Lernen

Optimierung und Generalisierung VL6

Vorhersage VL4 VL5 Neuronale Netze

Modellierung VL3 Supervised Unsupervised

VL7

VL8

### Historie und Philosophie

VL1

Geschichte

Personen

KI und Gesellschaft

Kritische Aspekte

### Mathematik

VL2

Lineare Algebra

Statistik

Logik

Numerik

Analysis

# KI-Landkarte

## Künstliche Intelligenz

### Modellierung und Schlussfolgerung

Variablen VL12 Inferenz

Logik Wissensrepräsentation

VL11

Zustände Suche VL13 MDPs

Reflex

### Anwendungen

Robotik VL14

Computer Vision VL9

Natürliche Sprache VL10

### Lernen

Optimierung und Generalisierung VL6

Vorhersage VL4 VL5 Neuronale Netze

Modellierung VL3 Supervised Unsupervised

VL7

VL8

### Historie und Philosophie

VL1

Geschichte

Personen

KI und Gesellschaft

Kritische Aspekte

### Mathematik

VL2

Lineare Algebra

Statistik

Logik

Numerik

Analysis

# KI-Landkarte

## Künstliche Intelligenz

### Modellierung und Schlussfolgerung

Variablen VL12 Inferenz

Logik Wissensrepräsentation

VL11

Zustände Suche VL13 MDPs

Reflex

### Anwendungen

Robotik VL14

Computer Vision VL9

Natürliche Sprache VL10

### Lernen

Optimierung und Generalisierung VL6

Vorhersage VL4 VL5 Neuronale Netze

Modellierung VL3 Supervised Unsupervised

VL7

VL8

### Historie und Philosophie

VL1

Geschichte

Personen

KI und Gesellschaft

Kritische Aspekte

### Mathematik

VL2

Lineare Algebra

Statistik

Logik

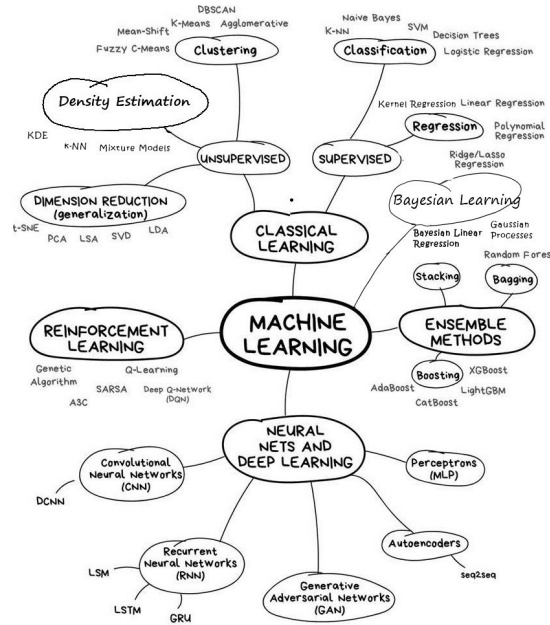
Numerik

Analysis

# Überblick

- 1) Lernende Algorithmen
- 2) Bewertung von lernenden Algorithmen
- 3) Klassifikation, logistische Regression

# Teil 1

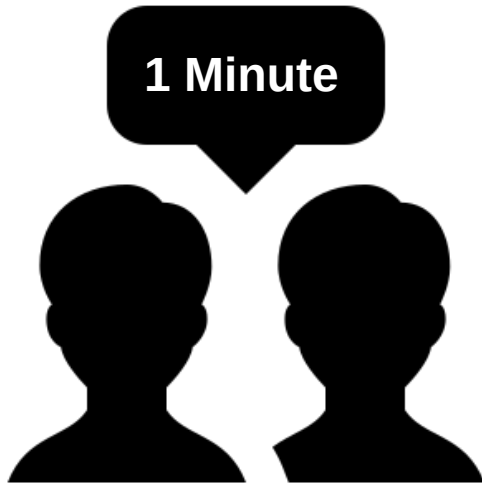


## 1) Lernende Algorithmen

2) Bewertung von lernenden Algorithmen

3) Klassifikation, logistische Regression

# Lernende Algorithmen - Definition



- Wie könnte man „Lernen“ definieren?
- Welche Begriffe müssen dazu wir definieren?

# Lernende Algorithmen - Definition

Mitchell 1997:

“A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .”

“Ein Computerprogramm kann als lernend bezeichnet werden, wenn es von Erfahrungen  $E$  bezüglich einer Klasse von Aufgaben aus  $T$  mit Leistungsmaß  $P$  lernt, und die Leistung in Aufgaben aus  $T$ , gemessen durch  $P$ , mit Erfahrung  $E$  zunimmt.

Was sind Beispiele für  **$T$  (tasks)**,  **$P$  (performance)** und  **$E$  (experience)**?

# Aufgaben (tasks), $T$

- Klassifikation
  - $k$  Kategorien von Input Vektoren:  $f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$
  - Beispiele: Objekterkennung in Bildern
- Regression
  - Vorhersage von Werten:  $f : \mathbb{R}^n \rightarrow \mathbb{R}$
  - Beispiele: Vorhersage von Mietpreisen
- Transkription
  - Unstrukturierter Input zu Text
  - Beispiele: Handschrift  $\rightarrow$  Text, Gesprochene Sprache  $\rightarrow$  Text
- Übersetzung
  - Symbole einer Sprache  $\rightarrow$  Symbole einer anderen Sprache
  - Beispiel: Natürliche Sprachen
- Ausreißer Erkennung
  - Kreditkartentransaktionen
- Synthese
  - Textgenerierung

# Leistungsmaß (performance), $P$

- Benötigt zur **quantitativen Analyse der Leistung** eines ML Algorithmus
- **Beispiele**
  - Klassifikation: Genauigkeit: Anteil der richtigen Vorhersagen
  - Regression: Maß für die Fehler der Vorhersagen
- **Allgemein: Fehler auf ungesehenen Daten**
  - Testset notwendig
  - Die Daten im Testset sollen aus der selben Verteilung kommen, wie die Daten auf die das Modell später angewendet werden soll
- Die negative Performance (Fehlermaß) wird oft als „Loss“ bezeichnet

# Erfahrung, $E$

- Typischerweise: **Datensatz** („dataset“)
  - Besteht aus  $m$  Datenpunkten, Samples
  - Oft: Jeder Datenpunkt ist ein Vektor mit  $n$  Einträgen
  - Jeder Eintrag wird als Attribut oder Feature bezeichnet
  - Design Matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$  : Ein Sample pro Zeile
- Je nach Typ von Daten: Einteilung der Lernverfahren in überwachte („supervised“) und unüberwachte („unsupervised“) Tasks

- **Supervised**: Jeder Datenpunkt hat ein Label oder einen Zielwert („target“)

$$p(\mathbf{y}|\mathbf{x})$$

- **Unsupervised**: Keine Labels oder Zielwerte

$$p(\mathbf{x})$$

# Beispiel (überwacht): Iris-Datensatz

Einer der ältesten Datensätze

150 Datenpunkte/Samples

3 Klassen, 4 Attribute/Features

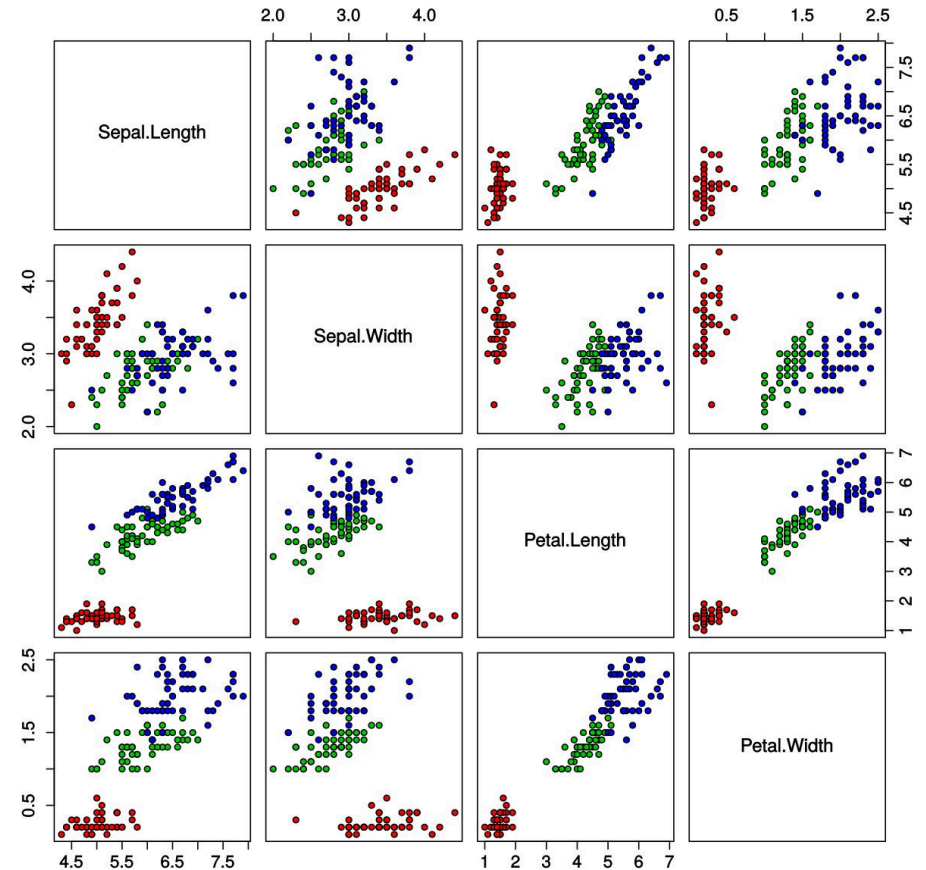
$$p(\mathbf{y}|\mathbf{x})$$



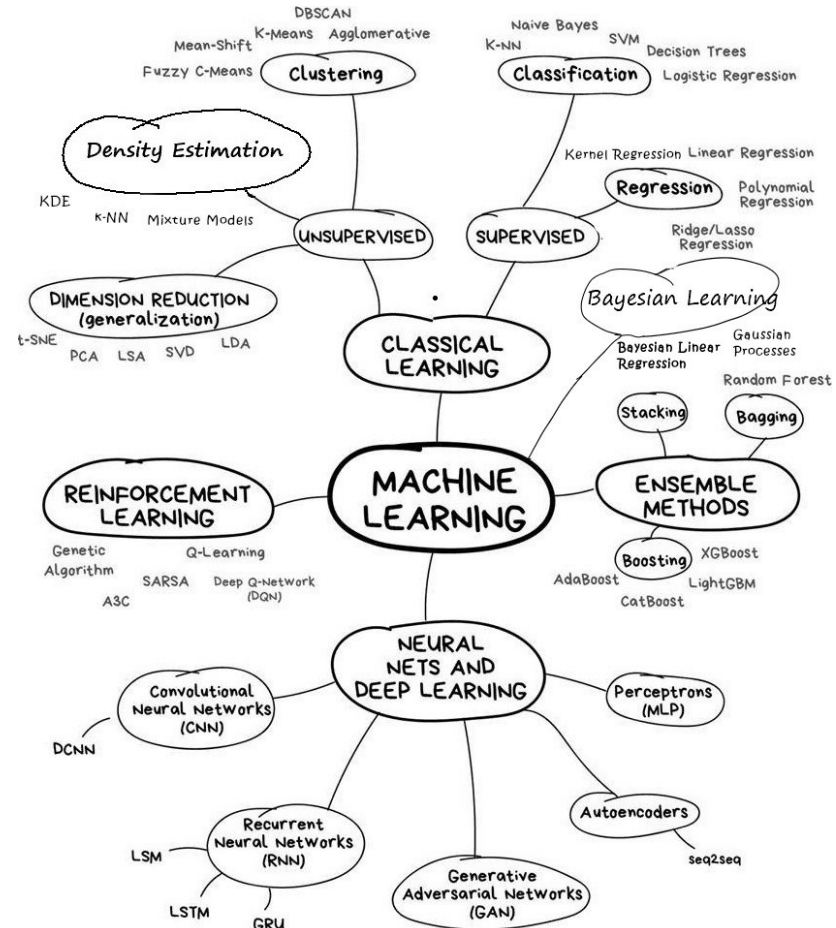
```
from sklearn.datasets import load_iris
```

```
iris = load_iris()
```

Iris Data (red=setosa,green=versicolor,blue=virginica)



# Machine Learning: Viele Begriffe und Methoden



# Gliederung nach **Lernmethoden** und Datentypen

- **Lernmethoden und Datentypen**

- **Supervised Learning**

*Daten mit Labels*

- Regression
    - Klassifikation

- **Unsupervised Learning**

*Daten ohne Labels*

- Clustering
    - Dimensionsreduktion
    - Density Estimation

- **Reinforcement Learning**

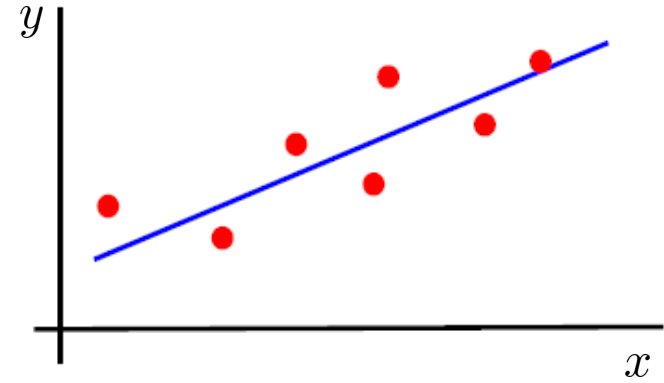
*Daten aus Interaktion*



# Supervised Learning

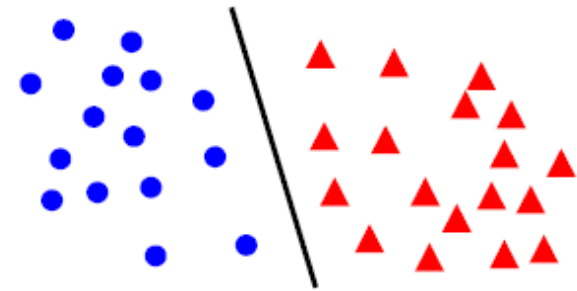
## 1. Regression

- Daten: Input  $x$  und Output („Labels“)  $y$
- Lernen einer kontinuierlichen Funktion  $y = f(x)$
- Beispiele: Lineare Regression



## 2. Klassifikation

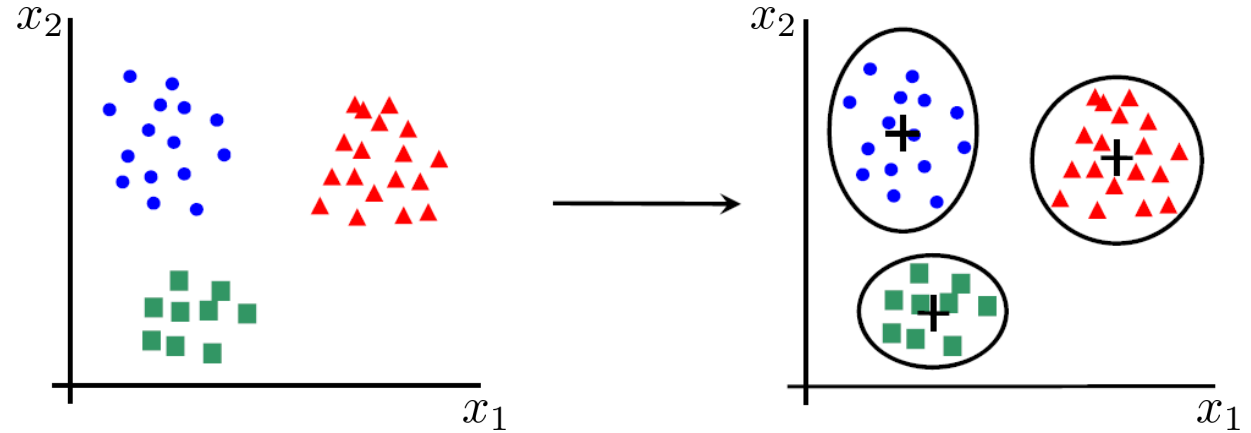
- Daten: Input  $x$  und Outputklassen („Labels“)  $c$
- Lernen der Klassenzuordnung  $c = f(x)$
- Beispiele: Handgeschriebene Zahlen



# Unsupervised Learning

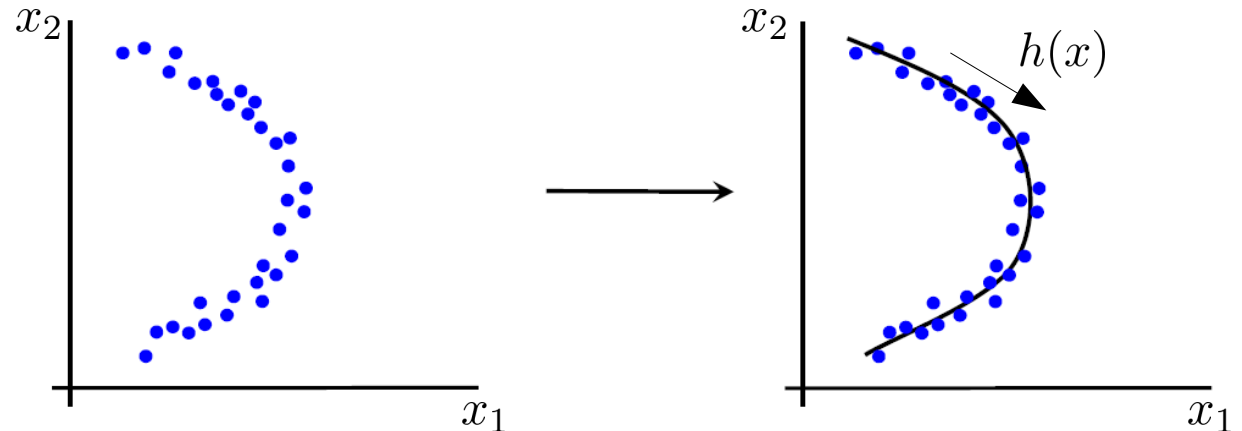
## 1. Clustering

- Daten:  $x$
- Lerne mögliche Klassen  $c$



## 2. Dimensionalitätsreduktion

- Daten: hochdimensionale  $x$
- Lerne niedrigdimensionale  $h$
- Lerne Funktion  $h(x)$  und  $x'(h)$



# Supervised vs. Unsupervised Learning

- Keine klare Unterscheidung
- Unüberwachtes Problem kann in  $n$  überwachte Probleme umformuliert werden

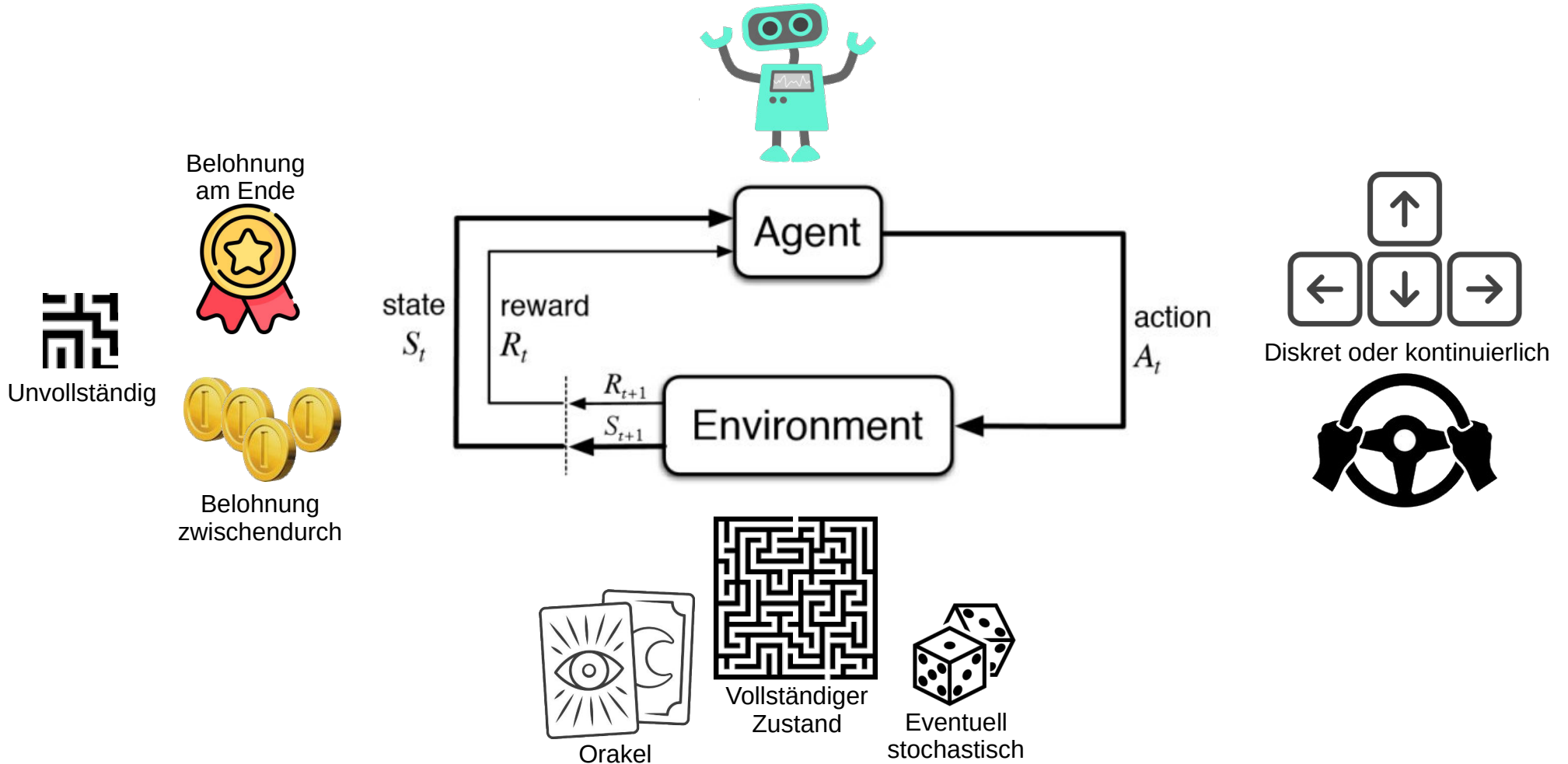
$$p(\mathbf{x}) = p(\mathbf{x}_0) p(\mathbf{x}_1|\mathbf{x}_0) p(\mathbf{x}_2|\mathbf{x}_0, \mathbf{x}_1) \dots = \prod_{i=1}^n p(\mathbf{x}_i|\mathbf{x}_1, \dots, \mathbf{x}_{i-1})$$

- Überwachtes Problem  $p(y|\mathbf{x})$  kann mit unüberwachtem Problem  $p(\mathbf{x}, y)$  gelöst werden

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}, y)}{\sum_{y'} p(\mathbf{x}, y')}$$

- Weitere Algorithmen können gar nicht eindeutig eingruppiert werden:
  - Bestärkendes Lernen (Reinforcement Learning) → Nächste Folie

# Reinforcement Learning



# Gliederung nach Modelltypen

- **Modelltypen**
  - **Symbolische Modelle**
    - Entscheidungsbäume
  - **Neuronale Netze, Deep Learning**
    - Multilayer Perceptrons
    - CNNs
    - RNNs
    - Autoencoder
    - GANs
  - **Kernel-Methoden**
    - Support Vektor Maschinen
    - Gaussian Processes



# Gliederung nach Modelltypen

- **Modelltypen**

- **Symbolische Modelle**

- Entscheidungsbäume

- **Neuronale Netze, Deep Learning**

- Multilayer Perceptrons
    - CNNs
    - RNNs
    - Autoencoder
    - GANs

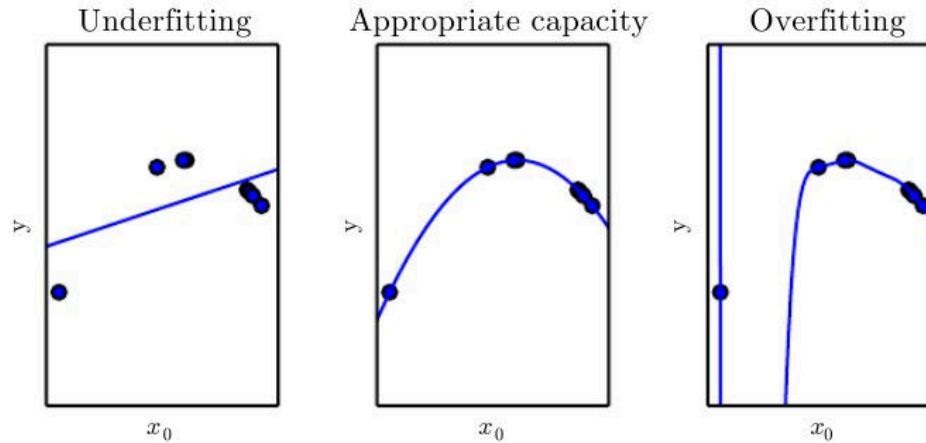
In dieser Vorlesung  
**Heute: Lineare Modelle**

- **Kernel-Methoden**

- Support Vektor Maschinen
    - Gaussian Processes



# Teil 2



1) Lernende Algorithmen

**2) Bewertung von lernenden Algorithmen**

3) Klassifikation, logistische Regression

# Kurze Erinnerung: Lineare Regression als ML Modell

- Gegeben:  $m$  Samples  $\mathbf{x} \in \mathbb{R}^n$  und  $m$  Zielwerte  $y \in \mathbb{R}$
- Modell:  $\hat{y} = \mathbf{w}^\top \mathbf{x}$  mit Parametervektor  $\mathbf{w} \in \mathbb{R}^n$
- Task  $T$ : Sage  $y$  für jedes  $\mathbf{x}$  voraus.
- Leistungsmaß  $P$ :

- Lege  $m_{\text{test}}$  Samples zurück: Testset  $\mathbf{X}^{(\text{test})}$  und  $\mathbf{y}^{(\text{test})}$

- Berechne Fehler auf dem Testset: 
$$\text{MSE}_{\text{test}} = \frac{1}{m} \sum_i \left( \hat{\mathbf{y}}^{(\text{test})} - \mathbf{y}^{(\text{test})} \right)_i^2$$

$$\text{MSE}_{\text{test}} = \frac{1}{m} \left\| \hat{\mathbf{y}}^{(\text{test})} - \mathbf{y}^{(\text{test})} \right\|_2^2$$

- Lösung durch Minimierung des Fehlers:  $\nabla_{\mathbf{w}} \text{MSE}_{\text{train}} = 0$

- Lösung: Normalengleichung:

$$\mathbf{w} = \left( \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \right)^{-1} \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})}$$

# Randnotiz: Bias in der linearen Regression

- Lineare Regression enthält im Allgemeinen auch einen Bias-Term:  $\hat{y} = \boldsymbol{w}^\top \boldsymbol{x} + b$
- Kann auch mit Normalgleichung gelöst werden:

$$\boldsymbol{x} = (x_1, \dots, x_n) \longrightarrow \boldsymbol{x} = (x_1, \dots, x_n, 1)$$

$$\boldsymbol{\omega} = (\omega_1, \dots, \omega_n) \longrightarrow \boldsymbol{\omega} = (\omega_1, \dots, \omega_n, b)$$

# Generalisierung

- Herausforderung des maschinellen Lernens: **Generalisierung**
- → Gute Voraussagen auch für ungesehene Daten (Fehler auf dem Testset)
  - **Generalisierungsfehler**
- Hauptunterschied zwischen ML und Optimierung
- Beispiel: Lineare Regression

- Training: Minimiere 
$$\text{MSE}_{\text{train}} = \frac{1}{m^{(\text{train})}} \left\| \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right\|_2^2$$

- Ziel: Minimiere 
$$\text{MSE}_{\text{test}} = \frac{1}{m^{(\text{test})}} \left\| \mathbf{X}^{(\text{test})} \mathbf{w} - \mathbf{y}^{(\text{test})} \right\|_2^2$$

- Lösung nur möglich, unter der **Annahme**, dass beide Datensätze aus der gleichen Verteilung gezogen werden

# Kapazität, Überfitten und Unterfitten

- Herangehensweise
  1. Minimiere den Fehler auf den Trainingsdaten
  2. Minimiere die Lücke zwischen Trainings und Testdaten
- Damit verbunden sind zwei grundlegende Konzepte des maschinellen Lernens:
  1. **Underfitting** → Großer Fehler auf Trainingsdaten
  2. **Overfitting** → Große Lücke zwischen Trainings- und Testfehler
- Kontrolle durch **Kapazität** eines Modells
  - Die Möglichkeit eines Modells, viele unterschiedliche Funktionen zu beschreiben
  - Niedrige Kapazität: Tendenz zum Unterfitten
  - Hohe Kapazität: Tendenz zum Überfitten

# Kapazität, Überfitten und Unterfitten

- Beispiel

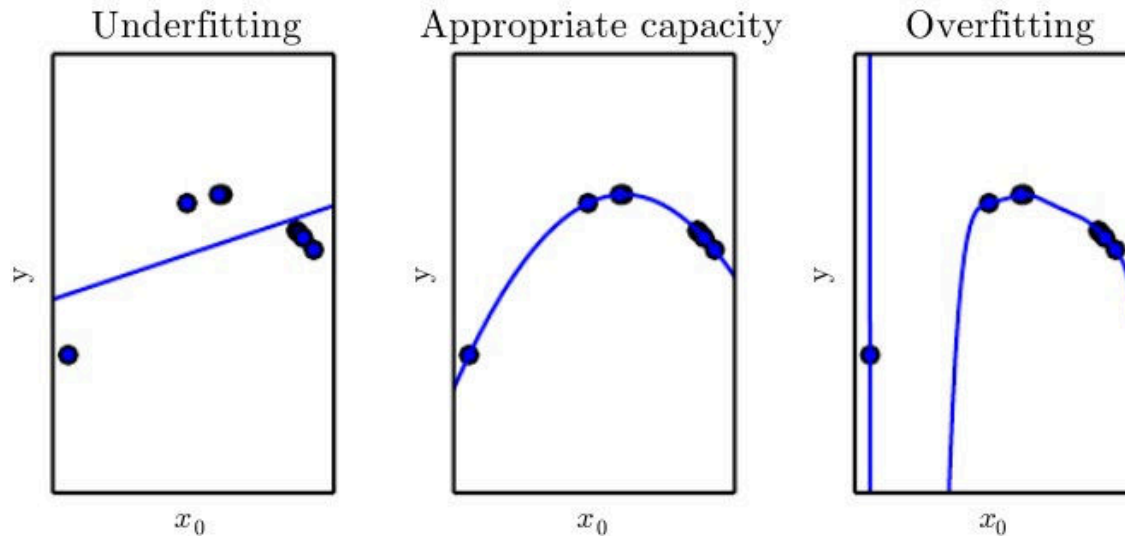
$$\hat{y} = b + wx$$

vs.

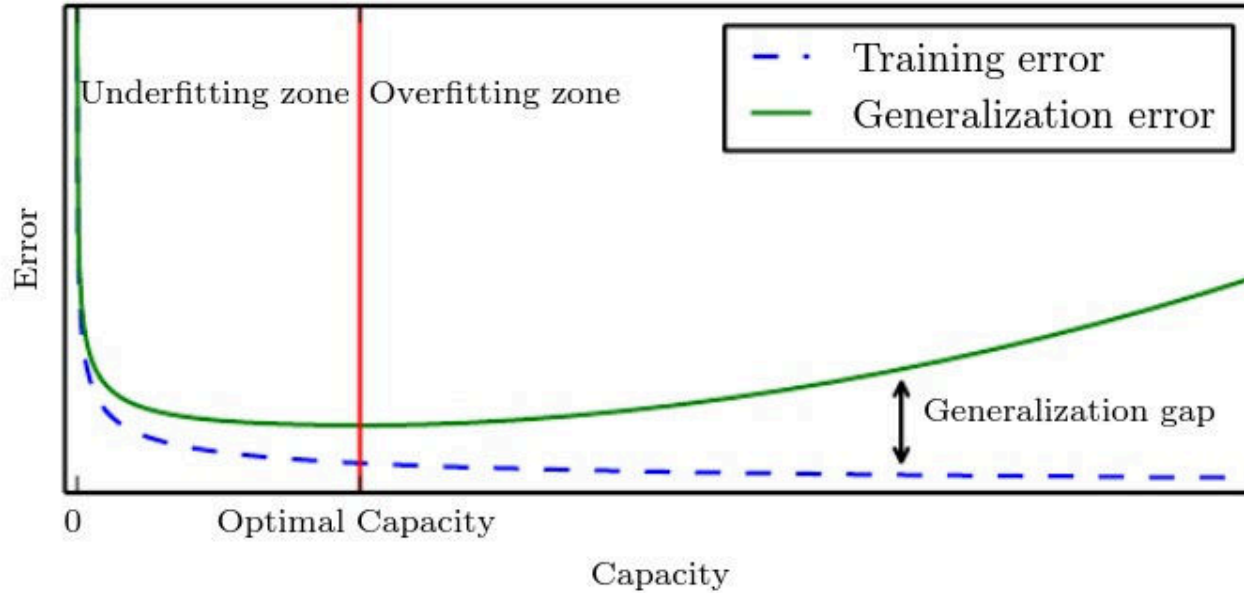
$$\hat{y} = b + w_1x + w_2x^2$$

vs.

$$\hat{y} = b + \sum_{i=1}^9 w_i x^i$$



# Kapazität, Überfitten und Unterfitten



# Begrenzung der Kapazität: Regularisierung

- Bisher: Kontrolle der Kapazität durch Änderung der zur Verfügung stehenden Funktionen
  - Benötigt Intuition und erhöht eventuell den Bias
- Alternativ: Modell muss sich selbst entscheiden
  - Bestrafung falls das Modell zu viele Funktionen benutzt
  - Quantisierung durch Gewichte („weight decay“)

$$\text{Loss}(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}$$

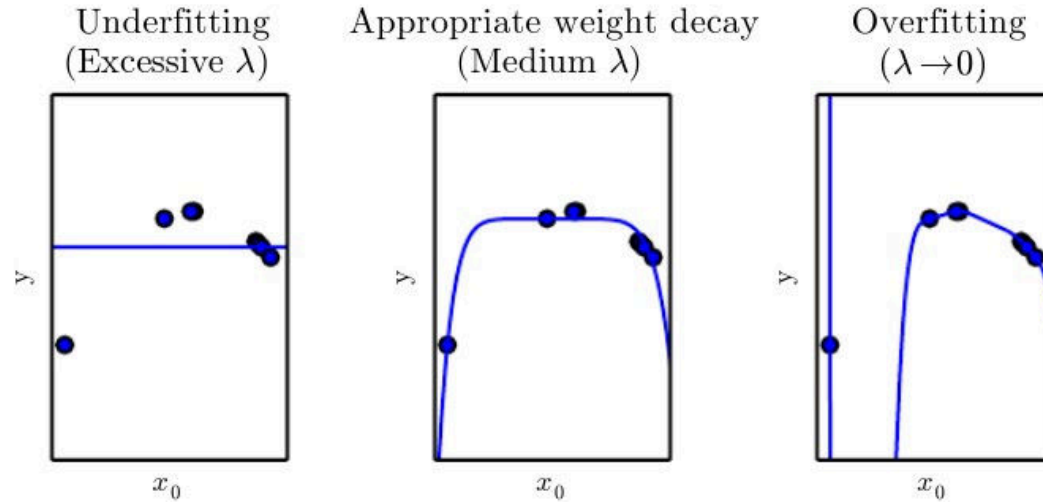
- Präferiert
  - Kleinere Gewichte
    - Dadurch weniger steile Funktionen („Steilheit = Ableitung = Gewichte“)
    - Glattere/einfacherer Funktionen

# Regularisierung

## Beispiel

- Datenverteilung: Quadratische Funktion
- Fit durch Polynom 9ten Grades
- Variation des Regularisierungsparameters  $\lambda$

$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}$$



# Warum überhaupt MSE? → Maximum-Likelihood-Methode

- **Annahme:** Gute Schätzfunktion (Parameter  $\theta$ ) durch Minimierung von  $\text{MSE} = \frac{1}{m} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2$
- **Allgemeiner**
  - Methode der maximalen Plausibilität (Maximum-Likelihood-Methode)
- Gegeben
  - $m$  Datenpunkte  $\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  verteilt mit unbekannter Verteilung  $p_{\text{data}}(\mathbf{x})$
- Gesucht
  - Parameter einer Wahrscheinlichkeitsfunktion  $p_{\text{model}}(\mathbf{x}; \theta)$

$$\begin{aligned}\theta_{\text{ML}} &= \arg \max_{\theta} p_{\text{model}}(\mathbb{X}; \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta)\end{aligned}$$

# Maximum-Likelihood-Methode

- Lösung von:  $\theta_{\text{ML}} = \arg \max_{\theta} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta)$

- Produkt durch Summe ersetzen  $\rightarrow$  Logarithmus:

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta)$$

- Teile durch  $m$ :

$$\theta_{\text{ML}} = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \theta)$$

- Dies entspricht einer Maximierung der Ähnlichkeit von empirischer Datenverteilung  $\hat{p}_{\text{data}}(\mathbf{x})$  und  $p_{\text{model}}(\mathbf{x}; \theta)$

# Maximum-Likelihood-Methode

- Für überwachtetes Lernen: Suche  $p_{\text{model}}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$

- Damit:  $\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbf{Y}|\mathbf{X}; \boldsymbol{\theta})$

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

- Regression mit fehlerbehafteten Daten

$$p_{\text{model}}(y|\mathbf{x}; \boldsymbol{\theta}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \boldsymbol{\theta}), \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y - \hat{y}(\mathbf{x}; \boldsymbol{\theta}))^2\right)$$

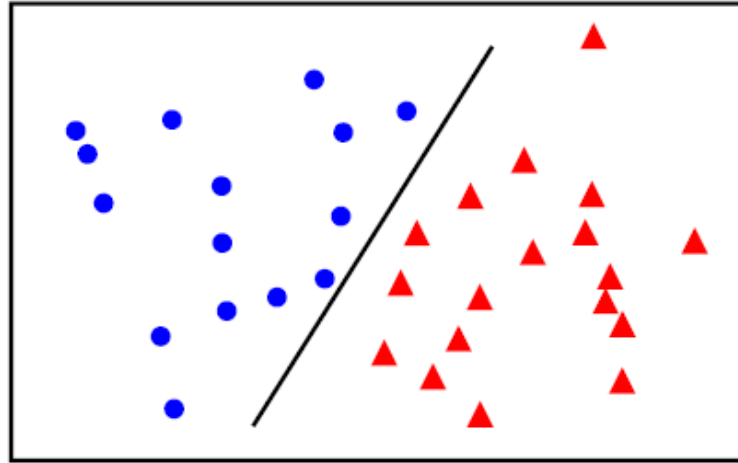
- Maximiere:  $\sum_{i=1}^m \log p_{\text{model}}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta})$

$$= -m \log \sigma - \frac{m}{2} \log(2\pi) - \sum_{i=1}^m \frac{\|y^{(i)} - \hat{y}(\mathbf{x}^{(i)}; \boldsymbol{\theta})\|^2}{2\sigma^2}$$

- Äquivalent zu Minimierung von  $\text{MSE}_{\text{train}} = \frac{1}{m} \sum_{i=1}^m \|y^{(i)} - \hat{y}(\mathbf{x}^{(i)}; \boldsymbol{\theta})\|^2$

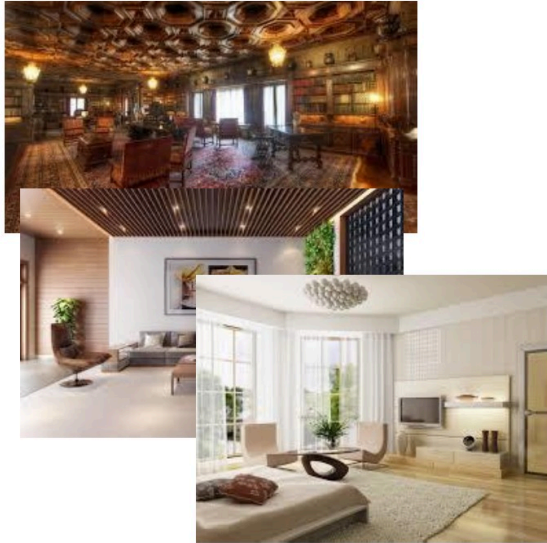


# Teil 3

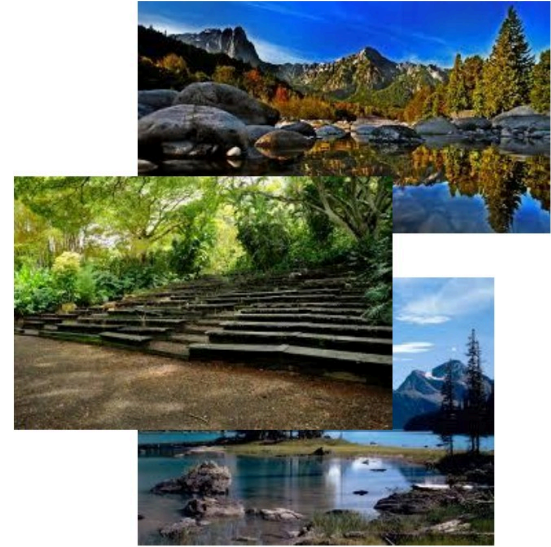


- 1) Lernende Algorithmen
- 2) Bewertung von lernenden Algorithmen
- 3) Klassifikation, logistische Regression**

# Beispiel: Klassifizierung von Bildern



Innen



Außen

(Bilddaten: Vorlesung 9)

# Beispiel: Spam Erkennung

	#"\$"	#"Mr."	#"sale"	...	Spam?
Email 1	2	1	1		Yes
Email 2	0	1	0		No
Email 3	1	1	1		Yes
...					
Email n	0	0	0		No
New email	0	0	1		??

# Klassifikation: Formale Definition

**Gegeben:** Datensatz  $\mathcal{D} = \{(\mathbf{x}_i, c_i)\}_{i=1\dots N}$   
mit Input Samples  $\mathbf{x}_i \in \mathbb{R}^n$  und Klassen  $c \in \{1 \dots K\}$

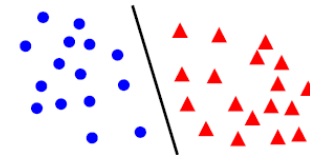
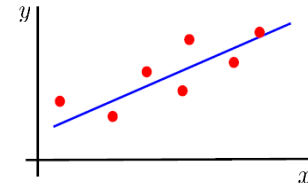
**Ziel:** Lernen einer Funktion  $c = f(\mathbf{x})$

**Nomenklatur:**  $K = 2$ : Binäre Klassifikation

$K > 2$ : Multiklassen Klassifikation

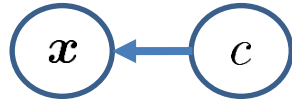
**Vergleich:** Regression: **Kontinuierlicher Output**

Klassifikation: **Diskreter Output**

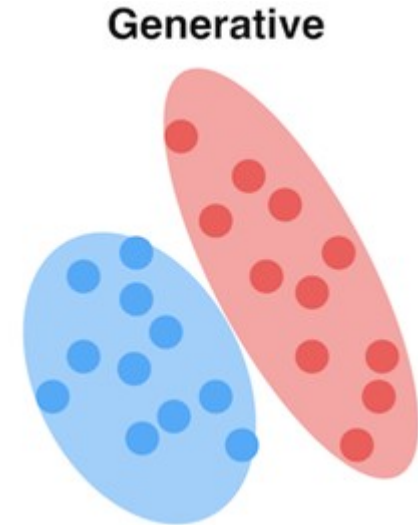


# Generative vs. diskriminative Modelle

## Generative Modelle

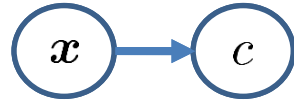


- Annahme:  $p(c)$
- Lernen eines Modells  $p(\mathbf{x} | c)$  aus den Daten
- Damit: Erzeugung von neuen Datenpunkten mit gegebenen Klassen
- Vorhersage der Klassen:  $p(c | \mathbf{x}) = \frac{p(\mathbf{x} | c)p(c)}{p(\mathbf{x})}$
- Das Lernen dieser Wahrscheinlichkeitsverteilungen ist sehr schwer!



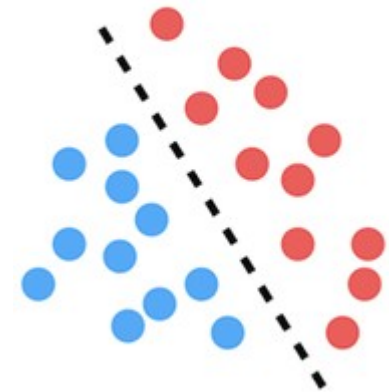
# Generative vs. diskriminative Modelle

## Diskriminative Modelle

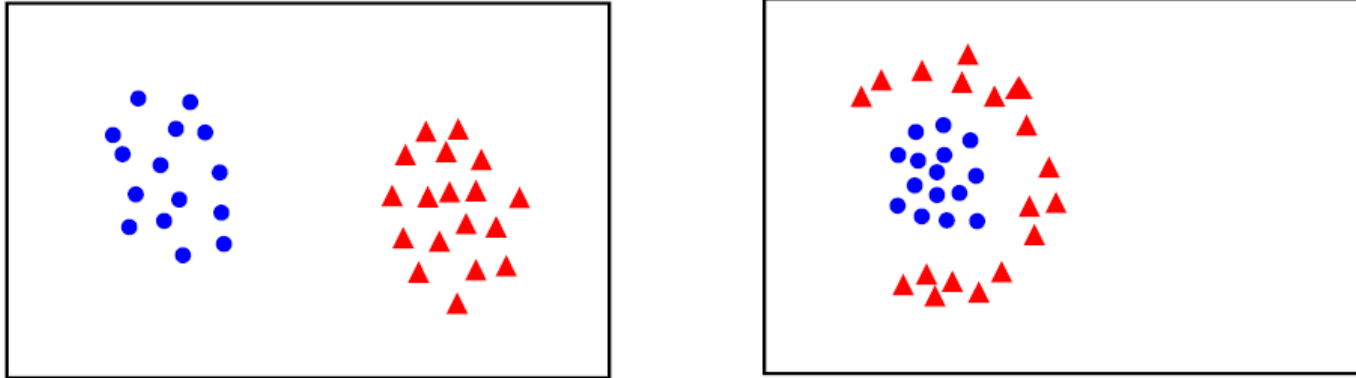


- Lernen einer Funktion  $f(x)$  oder einer Wahrscheinlichkeitsverteilung  $p(c | x)$
  - Die Funktion kann verschiedene  $x$  unterscheiden
  - Modellierung muss die Punkte an der Grenze berücksichtigen
  - Normalerweise sehr viel einfacher als Generation
- Wir fokussieren uns heute auf diskriminative Modelle

Discriminative



# Binäre Klassifikation



## Diskriminative Modelle

- Daten:  $(\mathbf{x}_i, c_i), i = 1 \dots m$  mit  $\mathbf{x}_i \in \mathbb{R}^n$  und  $c_i \in \{0, 1\}$

- Lerne einen Klassifikator:

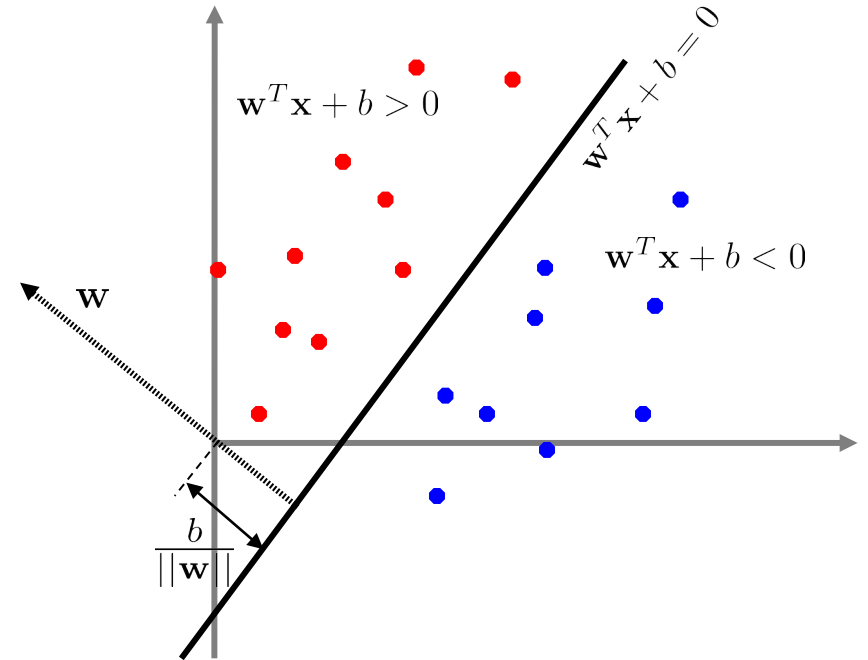
$$f(\mathbf{x}_i) = \begin{cases} > 0, & \text{if } c_i = 1 \\ < 0, & \text{if } c_i = 0 \end{cases}$$

# Lineare Klassifikation

## Linearer Klassifikator

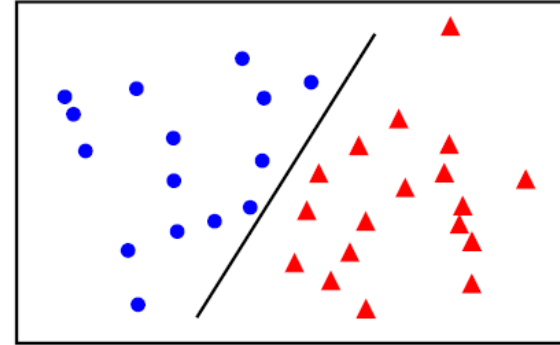
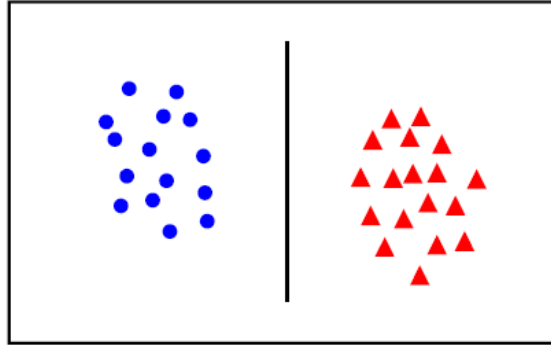
$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

- In 2D: Linie
  - $\mathbf{w}$  ist orthogonal zur Linie
  - $b$  ist der y-Achsen-Abschnitt („bias“)
- In 3D: Ebene
- In nD: Hyperebene

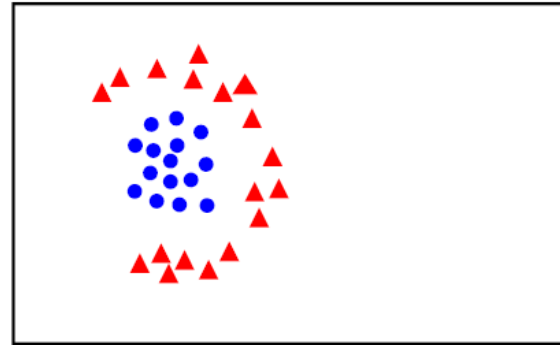
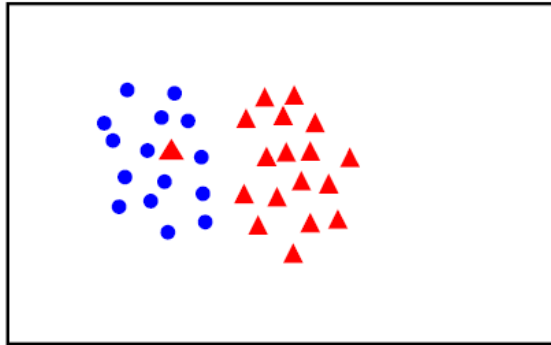


# Lineare Separarabilität

Linear separierbar



Nicht linear separierbar



# Lineare Klassifizierung: Erster Versuch

**Vorhersage:**  $y = \text{step}(f(\mathbf{x})) = \text{step}(\mathbf{w}^T \mathbf{x} + b)$

- Falls  $f(\mathbf{x}) < 0$  : Klasse 0
- Falls  $f(\mathbf{x}) \geq 0$  : Klasse 1

## Optimierung

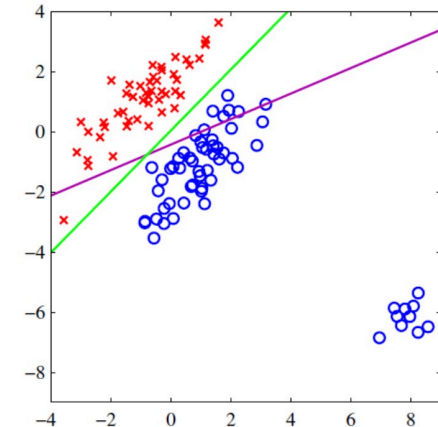
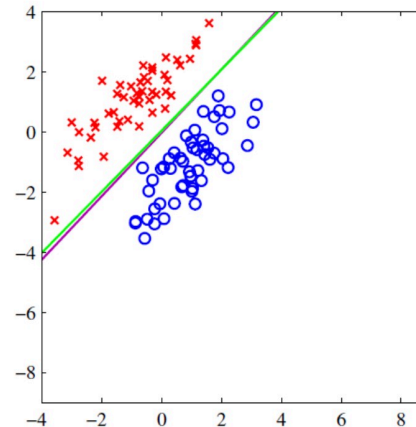
- Finde  $\mathbf{w}$  , sodass die Zahl der falsch klassifizierten Punkte minimal ist
- Quantifizierung:  $L_0(\mathbf{w}) = \sum \mathbb{I}(\text{step}(\mathbf{w}^T \mathbf{x} + b) \neq c_i)$
- Sehr schwer zu optimieren! NP-schwer!

# Lineare Klassifizierung: Zweiter Versuch

**Können wir die gleiche Loss-Funktion wie für Regression benutzen?**

$$L_{\text{reg}}(\mathbf{w}) = \sum_i (f(\mathbf{x}_i) - c_i)^2$$

- Minimierung der Loss-Funktion: Einfach!
- Aber: Die Labels sind nur 0 oder 1, der Output von  $f(\mathbf{x})$  ist kontinuierlich
- → Nicht robust gegenüber Ausreißern



# Lineare Klassifizierung: Dritter Versuch

## Können wir $f(x)$ einschränken?

- Benutze Sigmoid Funktion / logistische Funktion:

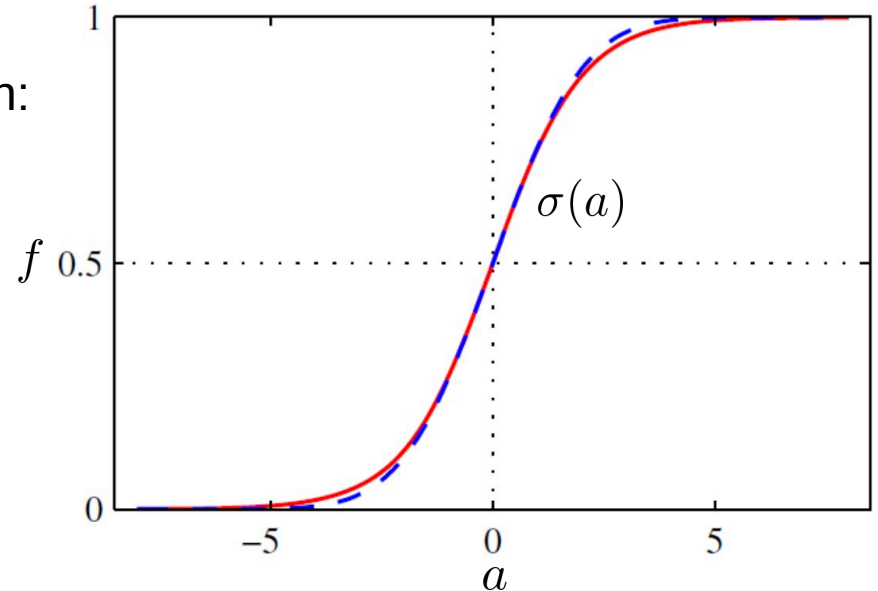
$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

- Output ist zwischen 0 und 1
- Keine Step-Funktion, sondern glatt

## Benutzung für Klassifikation

- Wende Sigmoid auf lineare Funktion an

$$L(\mathbf{w}) = \sum_i (\sigma(f(\mathbf{x}_i)) - c_i)^2 = \sum_i (\sigma(\mathbf{w}^T \mathbf{x}_i + b) - c_i)^2$$

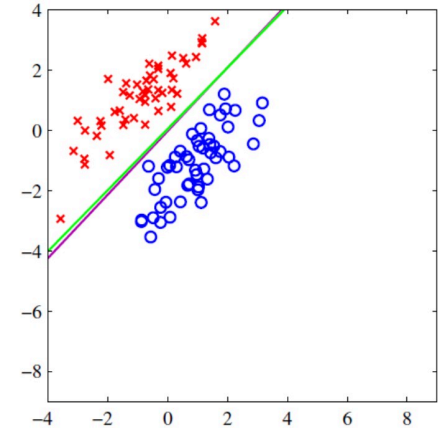


# Besser: Probabilistische Betrachtung

Interpretation von  $\sigma(f(x))$  als Wahrscheinlichkeitsfunktion:

$$p(c = 1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

$$p(c = 0 | \mathbf{x}) = 1 - \sigma(\mathbf{w}^T \mathbf{x} + b)$$



Trick um dies in einer Formel zusammenzufassen: Label/Klasse  $c$  im Exponent

$$p(c | \mathbf{x}) = p(c = 1 | \mathbf{x})^c p(c = 0 | \mathbf{x})^{1-c} = \sigma(\mathbf{w}^T \mathbf{x} + b)^c (1 - \sigma(\mathbf{w}^T \mathbf{x} + b))^{1-c}$$

(„Exponential-Trick“)

# Log-Likelihood

Maximierung der Wahrscheinlichkeit  $\rightarrow$  Maximierung des Logarithmus

$$\sum_i p(c_i | \mathbf{x}_i) \rightarrow \sum_i \log p(c_i | \mathbf{x}_i)$$

**Wir können direkt die log-likelihood optimieren (maximieren)**

$$\log \text{lik}(\tilde{\mathbf{w}}, D) = \sum_i \log p(c_i | \mathbf{x}_i) = \sum_i \log \left( p(c = 1 | \mathbf{x}_i)^{c_i} p(c = 0 | \mathbf{x}_i)^{1-c_i} \right)$$

$$= \sum_i c_i \log p(c = 1 | \mathbf{x}_i) + (1 - c_i) \log p(c = 0 | \mathbf{x}_i)$$

$$\tilde{\mathbf{x}}_i = [1, x_{i,0}, x_{i,1}, x_{i,2}, \dots, x_{i,n}]$$

$$\tilde{\mathbf{w}}_i = [b, w_0, w_1, w_2, \dots, w_n]$$

$$= \sum_i c_i \log \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) + (1 - c_i) \log (1 - \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i))$$

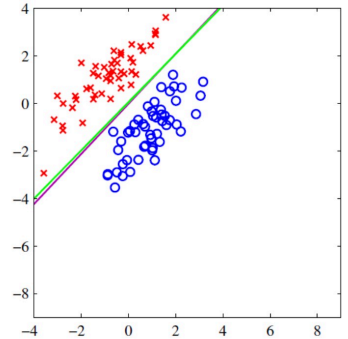
$\rightarrow$  Negative log-likelihood wird oft als „**cross-entropy loss**“ bezeichnet

# Logistische Regression

## Optimierung der log-likelihood der Sigmoid-Funktion: „logistische Regression“

$$\operatorname{argmax}_{\tilde{\mathbf{w}}} \log \operatorname{lik}(\tilde{\mathbf{w}}, D) = \operatorname{argmax}_{\tilde{\mathbf{w}}} \sum_i c_i \log \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) + (1 - c_i) \log (1 - \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i))$$

- ... obwohl wir es für Klassifikation (nicht Regression) benutzen
- Die Loss-Funktion ist konvex → Nur ein Optimum
- Es gibt keine geschlossene Lösung (anders als bei linearer Regression)



Wie finden wir das Maximum? → **Gradient descent** (Gradientenabstiegsverfahren)

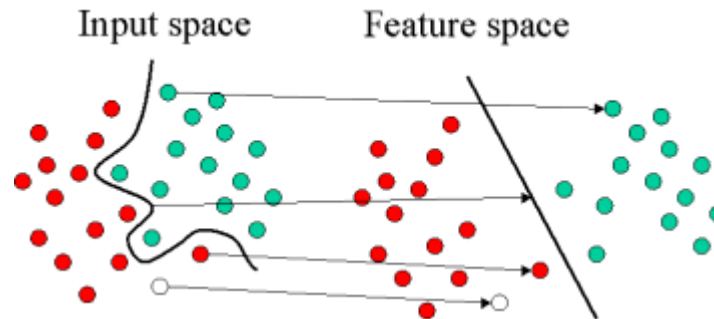
→ **Nächste Woche**

# Verallgemeinerte logistische Regression

## Linearer Klassifikator mit nicht-linearen Features $x_i \rightarrow \phi(x_i)$

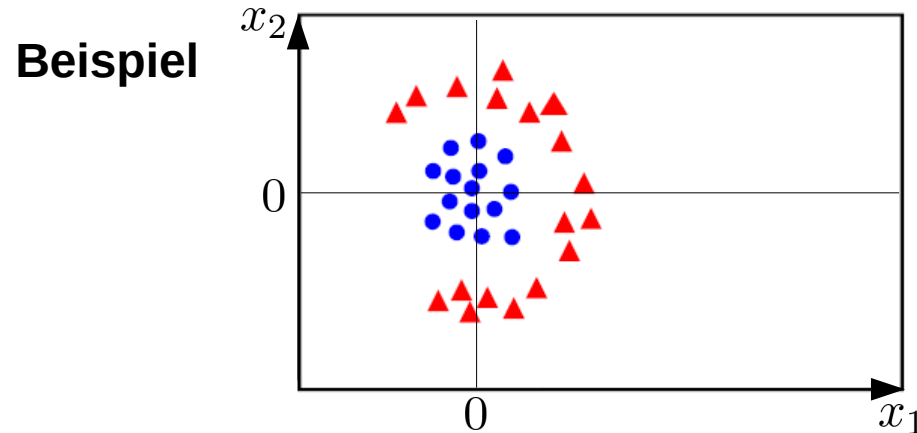
$$\operatorname{argmax}_{\mathbf{w}} \log \operatorname{lik}(\mathbf{w}, D) = \operatorname{argmax}_{\mathbf{w}} \sum_i c_i \log \sigma(\mathbf{w}^T \phi(\mathbf{x}_i)) + (1 - c_i) \log (1 - \sigma(\mathbf{w}^T \phi(\mathbf{x}_i)))$$

**Ziel:** Transformation der Punkte im Input-Raum (nicht linear separabel) zu Punkten im Feature-Raum die linear separabel sind



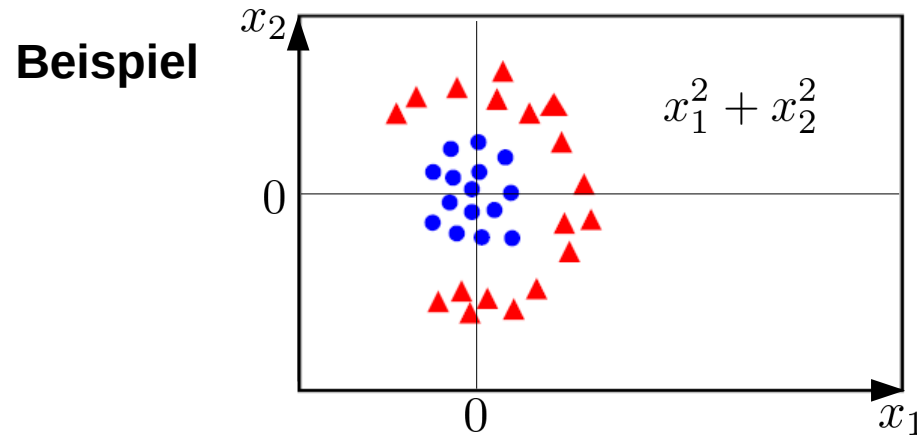
# Verallgemeinerte logistische Regression

**Ziel:** Transformation der Punkte im Input-Raum (nicht linear separabel) zu Punkten im Feature-Raum die linear separabel sind

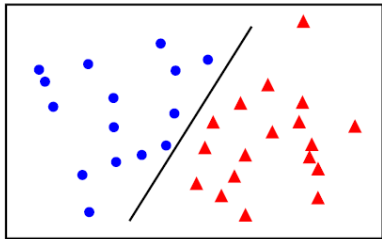


# Verallgemeinerte logistische Regression

**Ziel:** Transformation der Punkte im Input-Raum (nicht linear separabel) zu Punkten im Feature-Raum die linear separabel sind



# Zusammenfassung: Logistische Regression



Lineare Separierbarkeit

Sigmoidfunktion und logistische Regression

Log Likelihood Optimierung, Cross Entropy Loss

# Zusammenfassung

- 1) Lernende Algorithmen
- 2) Bewertung von lernenden Algorithmen
- 3) Klassifikation, logistische Regression

→ Nächste Woche: **Nicht-lineare Probleme, neuronale Netze!**