



Fakultät Ingenieurwissenschaften  
Labor für Kooperative, automatisierte Verkehrssysteme (KAV)

# **Verteidigungsmaßnahmen gegen Modell-Inversionsangriffe**

**Bachelorarbeit**

**von**

**Hannes Weber**

Aschaffenburg, 28. November 2023

Autor:

Hannes Weber

Am Lindenbrunnen 17

D-97846 Partenstein

Matrikel-Nr.: 2220472

Studiengang Software Design (Bachelor)

Prüfer:

Prof. Dr.-Ing. Konrad Doll

Kooperative, automatisierte Verkehrssysteme (KAV)

Zweitprüfer:

Prof. Dr.-Ing. Ulrich Brunsmann



Technische Hochschule Aschaffenburg  
Fakultät Ingenieurwissenschaften  
Würzburger Straße 45  
D-63743 Aschaffenburg

# Ehrenwörtliche Erklärung

Hannes Weber

Am Lindenbrunnen 17  
D-97846 Partenstein

Hiermit erkläre ich, dass ich die von mir vorgelegte Arbeit mit dem Thema „*Verteidigungsmaßnahmen gegen Modell-Inversionsangriffe*“ selbstständig verfasst habe, dass ich die verwendeten Quellen, Internet-Quellen und Hilfsmittel vollständig angegeben habe und dass ich die Stellen der Arbeit – einschließlich Tabellen und Bildern –, die anderen Werken oder dem Internet im Wortlaut oder dem Sinn nach entnommen sind, unter Angabe der Quelle als Entlehnung kenntlich gemacht habe.

Aschaffenburg, den 28. November 2023

---

Hannes Weber

# Danksagung

Hiermit möchte ich mich ausdrücklich bei Herrn Prof. Dr. Konrad Doll und Herrn Prof. Dr. Ulrich Brunsmann für die Unterstützung während meiner Arbeit bedanken. Des Weiteren ...

# Inhaltsverzeichnis

<b>Abkürzungen</b>	<b>1</b>
<b>1 Einleitung</b>	<b>2</b>
1.1 Motivation . . . . .	2
1.2 Aufgabenstellung . . . . .	3
1.3 Aufbau der Arbeit . . . . .	4
<b>2 Grundlagen</b>	<b>5</b>
2.1 Maschinelles Lernen . . . . .	5
2.1.1 Überwachtes Lernen . . . . .	6
2.1.2 Unüberwachtes Lernen . . . . .	7
2.1.3 Bestärkendes Lernen . . . . .	8
2.2 Bilderkennung /-klassifikation . . . . .	8
2.3 Angriffsmöglichkeiten auf Neuronale Netzwerke . . . . .	8
2.3.1 Modell-Inversionsangriffe . . . . .	8
<b>3 Stand der Technik</b>	<b>9</b>
3.1 Forschungsergebnisse . . . . .	9
<b>4 Implementierung</b>	<b>10</b>
4.1 Funktionalität des Codes . . . . .	10
4.2 Code-Beschreibung . . . . .	10
<b>5 Ergebnisse</b>	<b>11</b>
5.1 Beobachtungen . . . . .	11
5.2 Rückschlüsse . . . . .	12
<b>6 Zusammenfassung und Ausblick</b>	<b>13</b>
<b>7 Anhang</b>	<b>14</b>
<b>Bilderverzeichnis</b>	<b>16</b>
<b>Tabellen</b>	<b>17</b>
<b>Literatur</b>	<b>18</b>

## Abkürzungen

ML	<i>Maschinelles Lernen</i>
CNN	<i>Convolutional Neural Network</i>
KNN	<i>Künstlich-Neuronales Netzwerk</i>
NN	<i>Neuronales Netzwerk, eng.: Neural Network</i>

# 1

## Kapitel 1

---

### Einleitung

#### 1.1 Motivation

Inmitten der fortschreitenden Digitalisierung greift man immer häufiger auf Lösungen zurück, die diverse Aufgaben in Arbeits- und Lebensbereich unterstützen, vereinfachen und sogar ersetzen. Von Medizin, in der beispielsweise Diagnoseverfahren zur frühzeitigen Erkennung von Krankheiten eingesetzt werden, bis zur Automobilindustrie, die schon jetzt auf selbstfahrende Kraftfahrzeuge setzt. Durch all die neuen Entwicklungen und Einsatzgebiete von künstlichen Intelligenzen gewinnt diese unaufhaltsam an Bedeutung.

Trotz des enormen Potenzials, das diese Technologien in ihren jeweiligen Anwendungsgebieten mit sich bringen, treten vermehrt Herausforderungen hinsichtlich Sicherheit und Robustheit der Systeme auf. Ein zentraler Aspekt ist die Notwendigkeit der Risikominimierung durch das Sicherstellen der Vertrauenswürdigkeit von KI-Systemen. Ein weiterer wichtiger Punkt in der Entwicklung dieser Systeme ist die Robustheit gegenüber Angriffen. KI-Systeme können anfällig für Manipulationen und gezielte Angriffe von außen oder durch Fehlkonfigurationen sein. Dabei spielt die Implementierung und Integration von bestimmten Sicherheitsmechanismen eine entscheidende Rolle, um die Integrität der Systeme zu gewährleisten. Insgesamt ist die Sicherheit im Bereich des maschinellen Lernens von zentraler Bedeutung, wodurch das Vertrauen der Nutzer gestärkt und die breite Integration dieser Technologien in verschiedener Bereichen vorangetrieben wird. Diese Arbeit legt den Fokus auf die Herausforderung der Verdeutlichung von Sicherheits- und Robustheitsimplementierungen in Bild-Klassifikationsmodellen. In einem globalen System, in dem immer mehr Aspekte des täglichen Lebens in die 'Hände' von KI-Systemen gegeben werden, ist es umso wichtiger, diese mit Anbetracht auf Sicherheit zu implementieren, wie auch zu überwachen. Eine Unsicherheit eines Systems kann hierbei schon zu Verletzungen der Privatsphäre einzelner Personen führen. Daher wird diese Arbeit gesondert Angriffsvektoren und Bedrohungen von Klassifikationsmodellen behandeln. Darüber hinaus werden innovative Ansätze zur Bekämpfung möglicher Angriffsvektoren und Schwachstellen aufgezeigt, die sowohl Sicherheit, als auch Robustheit von neuronalen Netzen erhöhen. Insgesamt soll diese Arbeit dazu beitragen, dem Leser ein grundlegendes Verständnis über Angriffe, deren Auswirkungen, wie auch Verteidigungsmaßnahmen zu liefern.

## 1.2 Aufgabenstellung

### Hintergrund

Da KI-Systeme immer mehr Anwendung finden, werden diese zudem häufiger Ziele von Cyber-Angriffen. Es gibt viele verschiedene Angriffsvektoren auf diese Systeme, zu denen auch das Erlangen von zugrundeliegenden Informationen über Daten des Trainingsprozesses von Neuronalen Netzen gehört. Da viele Modelle mit sensiblen Daten trainiert werden, ist es wichtig, dass Daten durch das Modell nicht an Angreifer übergeben werden. Der 'EU-AI Act' enthält unter anderem Anforderungen an KI-Systeme in Bezug auf Robustheit und Cybersicherheit. Daher bringt dieser nicht nur die Verantwortung sichere Modelle zu trainieren mit sich, sondern auch die Sicherstellung, dass genutzte Trainingsdaten privat gehalten werden. Um Modelle während der Bereitstellung abzusichern, müssen die verschiedenen Angriffsvektoren und die entsprechenden Abwehrmaßnahmen bekannt sein.

### Ziel der Arbeit

Während der Thesis sollen folgende Fragen beantwortet werden:

- Welche Angriffsmöglichkeiten gibt es, um Daten von deployten Modellen zu extrahieren?
- Wie kann man sich gegen diese Attacks schützen?
- Showcase zu Verteidigungstechniken und deren Effektivität gegenüber Inversions-Angriffen.

### Methodischer Ansatz

- Onboarding
- Recherche über verschiedene Attacks und Verteidigungsmöglichkeiten
- Vergleich von verschiedenen Angriffen auf unterschiedliche Modell-Architekturen
- Implementierung eines Showcases für mindestens einen Angriff auf mindestens eine Modell-Architektur:
  - Zeigen, wie ein solcher Angriff funktioniert.
  - Kann man sich gegen einen solchen Angriff verteidigen?
  - Wie effektiv sind die Verteidigungsstrategien?



## 1.3 Aufbau der Arbeit

Die Arbeit unterteilt sich in drei Hauptbestandteile. Dazu gehört zum einen der SStand der Technik", worin aktuelle Technologien und Modelle dargestellt werden. Hauptsächlich wurde der Wissenstand aus anderen Papern entnommen und spiegelt somit den "ResearchingTeil der Arbeit wieder. ....

Eine weitere Hauptkomponente der Arbeit ist die Implementierung eines Modells und dessen Angriffsvektoren. Mit der Implementierung der möglichen Verteidigungsmaßnahmen der Schwachstellen wird die Kehrseite aufgezeigt, mit der eine mögliche Absicherung stattfinden kann. ...

Durch die Darstellung und Auswertung der herausgefundenen Ergebnisse, geht die Arbeit zu Ende. Dabei werden die wichtigsten Erkenntnisse aus den Implementierungen und Researching-Tasks zusammengefasst, und anschaulich dargestellt. Dies soll den Vorteil mit sich bringen, dem Leser einen Überblick des neu erworbenen und eine finale Aussicht zu bieten. ...

Der Aufbau der Arbeit lehnt sich an die Struktur von **Bar-Shalom** an.

# 2

## Kapitel 2

---

### Grundlagen

#### 2.1 Maschinelles Lernen

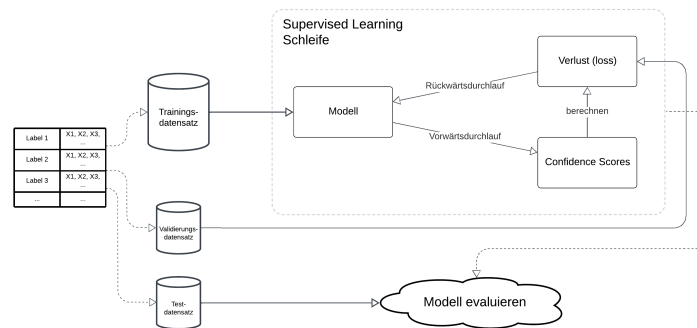
Ein Teilgebiet der künstlichen Intelligenzen ist der Bereich des Maschinellen Lernens, welcher die Verarbeitung von Daten adressiert, um beispielsweise Vorhersagen in Situationen auf Basis von bestimmten Informationen zu treffen. Im Jahr 1959 wurde dieser von Arthur Samuel geprägt, der die Herausforderung annahm, ein Schachspiel mit Hilfe einer Maschine zu lösen. (Joshi 2020, S. 4) Das besondere hierbei ist das Erlernen des richtigen Lösungswegs, welcher nicht durch bestimmte Bedingungen gesteuert, sondern auf Basis verschiedener Informationen erlernt wird, die über diverse Quellen beigefügt werden. Maschinelles Lernen findet in vielen Bereichen Anwendung, von Spracherkennung, Bilderkennung bis hin zu selbstfahrenden Kraftfahrzeugen. Diese Technologie wird zunehmend wichtiger in unserer immer vernetzteren Welt und trägt dazu bei, komplexer werdende Probleme zu lösen und immer innovativere Lösungen in verschiedenen Branchen einzuführen.

Im Wesentlichen geht es um die Entwicklung von Algorithmen, die dem Computer das Erkennen von Mustern und Zusammenhängen in Daten ermöglichen, wodurch ohne menschliche Intervention Aufgaben erledigt werden können. Daher trainiert man Algorithmen auf Basis großer Datenbestände, um die erlernten Fähigkeiten auf neu erhobene Daten anzuwenden. Das Erlernen des nötigen Wissens basiert hierbei auf drei wesentlichen Faktoren, die die Qualität des Trainings beeinflussen. Neben den Informationen, die aus Daten gewonnen werden, nutzt man eine Kennzahl, die einen Vergleich zwischen dem aktuellen und idealen Verhalten herstellt, um mit dem dritten Faktor - einem Rückkopplungsmechanismus - das Programm anzuleiten eine verbesserte Leistung in Folgeergebnissen zu erzielen. (Joshi 2020, S. 4)

Um ein Modell trainieren zu können, wird ein geeignetes künstliches Lernverfahren ausgewählt, das darauf abzielt, die Ausgaben eines künstlichen Systems in Bezug auf bestimmte Systemeingaben im Laufe des Lernprozesses zu optimieren. Unterschieden wird hauptsächlich in der Art und Weise, wie die 'Kritik' präsentiert wird, die zur Verbesserung des Verhaltens der jeweiligen künstlichen Systeme führen soll. Ein Training basiert auf sogenannten Hyperparametern, die Rahmenbedingungen des Prozesses festlegen. Dies kann zum Beispiel die Anzahl der Epochen, die Lern-Rate oder auch die Batch-Größe sein. Im Folgenden werden drei prägende Lernverfahren näher beschrieben.

### 2.1.1 Überwachtes Lernen

Diese Art des Lernens (eng.: supervised learning) stellt einen wesentlichen Bereich des maschinellen Lernens dar, um aus Informationen, bestehend aus Datenpunkten  $X = x_1, x_2, \dots, x_n$  mit einem zugehörigen Label aus  $Y = y_1, y_2, \dots, y_n$ , das nötige Wissen für bestimmte Anwendungsfälle zu erlangen. Anwendung findet dieses Verfahren in Systemen, die nicht-gelabelte Daten verarbeiten, um darauf basierend eine Vorhersage über die Zugehörigkeit abzugeben. Mögliche Arten von neuronalen Netzwerken, die über diese Methode trainiert wurden, sind beispielsweise Bildklassifikatoren. Dabei ist es die Aufgabe des Modells  $M$ , basierend auf einem bestimmten Input  $I$  mit Hilfe erlernter Muster und Zusammenhänge der Trainingsdaten eine Vorhersage über die Zugehörigkeit von  $I$  zu treffen.



**Bild 2.1:** Visualisierung grundlegender Prozesse während des Ablaufs eines überwachten Lernens

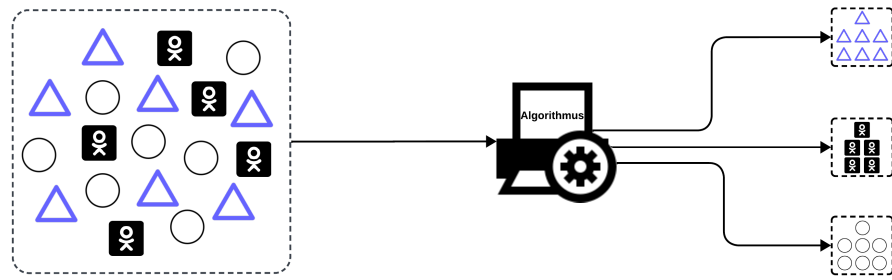
Der *Ablauf* des Trainings lässt sich wie folgt beschreiben: nach Beenden von Datensammlung, Datenvorverarbeitung und Datensatztrennung in einen Trainings-, Validierungs- und Testdatensatz beginnt der eigentliche Teil des Trainings, wobei die Trainings- und Validierungsdaten für die Aktualisierung der Modellparameter verwendet werden. Die Optimierung der Parameter wird auf Basis der Verlustfunktion mit Hilfe eines Rückwärtsthroughlaufs (eng.: Backward Propagation) durchgeführt, die mit der Differenz zwischen Vorhersagen des Modells aus dem Vorwärtsthroughlauf (eng.: Forward Propagation) und den realen Labels berechnet wird.

Die *Qualität* des Trainings wird meist anhand der Genauigkeit von Vorhersagen bestimmt, welche ausdrückt, wie gut die Klassifikation von im Training noch nicht gesehenen Daten funktioniert.

$$\text{Accuracy} = \frac{\text{Anzahl korrekter Vorhersagen}}{\text{Gesamtanzahl der Vorhersagen}}$$

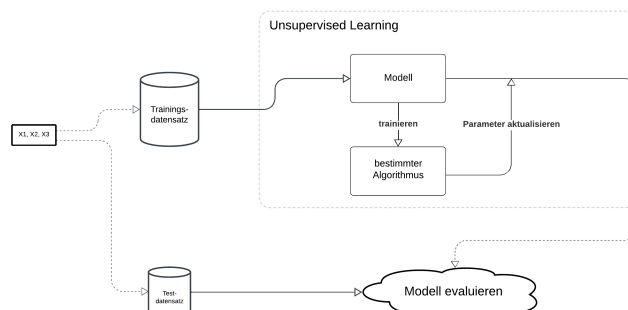
Weitere Metriken können die Präzision (eng.: Precision), der Recall oder auch der F1-Score sein. Neben Messungen der Modell-Genauigkeit kann man die Qualität des trainierten Modells auch anhand der Verlustfunktion (eng.: loss-function) oder der 'Area Under the Curve' (AUC), welche die Fläche unter der 'Receiver Operating Characteristic' (ROC) – dem Verhältnis zwischen der True- und False-Positive-Rate – messen.

## 2.1.2 Unüberwachtes Lernen



**Bild 2.2:** Resultat eines auf unüberwachtem Lernen basierenden Clusterings

Unüberwachtes Lernen (eng.: unsupervised learning) stellt einen anderen bedeutenden Bereich im maschinellen Lernen dar, der sich von überwachtem Lernen (2.1.1) dahingegen unterscheidet, dass keine expliziten Labels für die Trainingsdaten bereitgestellt werden. Diese Methode wird angewendet, wenn das Ziel darin besteht, Muster, Strukturen oder Zusammenhänge in den Daten zu entdecken, ohne dabei Kategorien oder Label vorzugeben, weshalb basierend auf Informationen eines Datensatzes, bestehend aus Datenpunkten  $X = x_1, x_2, \dots, x_n$  ohne dazugehöriges Label, trainiert wird. Das Modell soll dabei auf natürliche Art und Weise Strukturen, Muster und Zusammenhänge innerhalb der Eingabedaten ohne vorherige Kenntnisse der Zielvariable erlernen. „Aufgabe ist es hier, passende Repräsentationen zu finden, die z. B. die Erkennung von Charakteristika in Datenmengen, Wiedererkennung von Ausnahmen oder die Erstellung von Prognosen ermöglichen.“ (Lorenz 2020, S. 5) Im Gegensatz zu Überwachtem Lernen (2.1.1) geht es nicht um das Zurechnen von Mustern in vorhandenen Kategorien, sondern um das Auffinden von Clustern in einer bestimmten Datenmenge  $X$ . Dabei werden Modellparameter nicht über eine Verlustfunktion, die durch die Differenz zwischen tatsächlichem und vorhergesagtem Label berechnet wird, dargestellt, sondern die Parameter werden für die repräsentation inhärenter Muster in den Daten angepasst. Das Aufteilen des Datensatzes in 3 unabhängige Teile ist hierbei nicht notwendig, da keine Kategorien/Label für die jeweiligen Datenpunkte vorhanden sind. Es ist nur für die Auswertung des Algorithmus sinnvoll einen zweiten Datensatz zu erstellen, um damit die Qualität zu bewerten.



**Bild 2.3:** Visualisierung grundlegender Prozesse während des Ablaufs eines unüberwachten Lernens

Der *Ablauf* des Trainings gestaltet sich wie folgt: Nach Abschluss von Datensammlung und Vorverarbeitung beginnt der eigentliche Prozess des unüberwachten Lernens. Im Gegensatz zum überwachten Lernen (2.1.1) gibt es hierbei keine vordefinierten Zielvariablen. Den Datenpunkten  $X = x_1, x_2, \dots, x_n$  ist also initial keine Kategorie beziehungsweise kein Label zugeordnet, da die Menge  $Y$  bis dato nicht existent ist. Im Gegensatz zu überwachtem Lernen nutzt man hier Algorithmen wie beispielsweise k-means Clustering, um eine bestimmte Operation auf den zu behandelnden Datensatz auszuführen. Mit Hilfe des Algorithmus wird ein bestimmtes Modell dahingehend trainiert, dass es einen bestimmten Datensatz zum Beispiel in bestimmte Kategorien unterteilen kann. Dafür analysiert der Algorithmus die jeweiligen Datenpunkte und versucht diese mit Hilfe von Zusammenhängen und verschiedenen extrahierten Merkmalen zu interpretieren. Anhand des dabei erlangten Wissens werden Modellparameter aktualisiert, wonach das unüberwachte Lernen beendet ist. Ein Ergebnis kann hierbei beispielsweise ein kategorisierter Datensatz sein.

### **2.1.3 Bestärkendes Lernen**

## **2.2 Bilderkennung /-klassifikation**

## **2.3 Angriffsmöglichkeiten auf Neuronale Netzwerke**

Hier soll geschildert werden, welche Angriffsmöglichkeiten es gibt. Im ersten Unterpunkt soll auf Modell-Inversionsangriffe eingegangen werden.

### **2.3.1 Modell-Inversionsangriffe**

Hier soll was über Modell-Inversionsangriffe stehen (Bsp. usw.) ...

#### **2.3.1.1 Angriffsziel**

Hier soll was über Ziel von diesem Angriff geschrieben werden (Bsp. usw.) ...

#### **2.3.1.2 Angriffsvektoren**

Hier soll geschrieben werden, wie der Angriff ausgeführt werden kann (Bsp. usw.) ...

#### **2.3.1.3 Verteidigungsstrategien**

Hier sollen mögliche Verteidigungsstrategien beleuchtet werden ...

# 3

Kapitel 3

---

## Stand der Technik

### 3.1 Forschungsergebnisse

# 4

## Kapitel 4

---

### Implementierung

#### 4.1 Funktionalität des Codes

Hier soll die Funktionalität des Codes beschrieben werden ...

```
1  import numpy as np
2
3  def incmatrix(genl1,genl2):
4      m = len(genl1)
5      n = len(genl2)
6      M = None #to become the incidence matrix
7      VT = np.zeros((n*m,1), int) #dummy variable
8
9      #compute the bitwise xor matrix
10     M1 = bitxormatrix(genl1)
11     M2 = np.triu(bitxormatrix(genl2),1)
12
13     for i in range(m-1):
14         for j in range(i+1, m):
15             [r,c] = np.where(M2 == M1[i,j])
16             for k in range(len(r)):
17                 VT[(i)*n + r[k]] = 1;
18                 VT[(i)*n + c[k]] = 1;
19                 VT[(j)*n + r[k]] = 1;
20                 VT[(j)*n + c[k]] = 1;
21
22             if M is None:
23                 M = np.copy(VT)
24             else:
25                 M = np.concatenate((M, VT), 1)
26
27             VT = np.zeros((n*m,1), int)
28
29     return M
```

Listing 4.1: Python example

#### 4.2 Code-Beschreibung

Hier soll beschrieben werden, wie der Code strukturiert ist, nach welchen Pattern gearbeitet wurde und was der Code oberflächlich macht...





steht ganz viel Text... Hier steht ganz viel Text... Hier steht ganz viel Text... Hier steht ganz viel Text... Hier steht ganz viel Text... Hier steht ganz viel Text... Hier steht ganz viel Text... Hier steht ganz viel Text... Hier steht ganz viel Text...

## **5.2 Rückschlüsse**

Hier sollen die Rückschlüsse stehen ...

# 6

## Kapitel 6

---

### Zusammenfassung und Ausblick

# 7

## Kapitel 7

---

### Anhang

Hier sind noch  $\text{\LaTeX}$ -Beispiele für die Verwendung:

...

Tabelle erzeugen:

**Tabelle 7.1:** Beispieltabelle

Spalte 1	Spalte 2	Spalte 3	Spalte 4
Inhalt 1	Inhalt 2	Inhalt 3	Inhalt 4
Inhalt 5	Inhalt 6	Inhalt 7	Inhalt 8

...

Quellen verlinken: Test **Bar-Shalom**. Test **Goldhammer**. Test **WHO-2004**.

...

Kapitel verlinken: Wie in Abschnitt 5 beschrieben wird ...

...

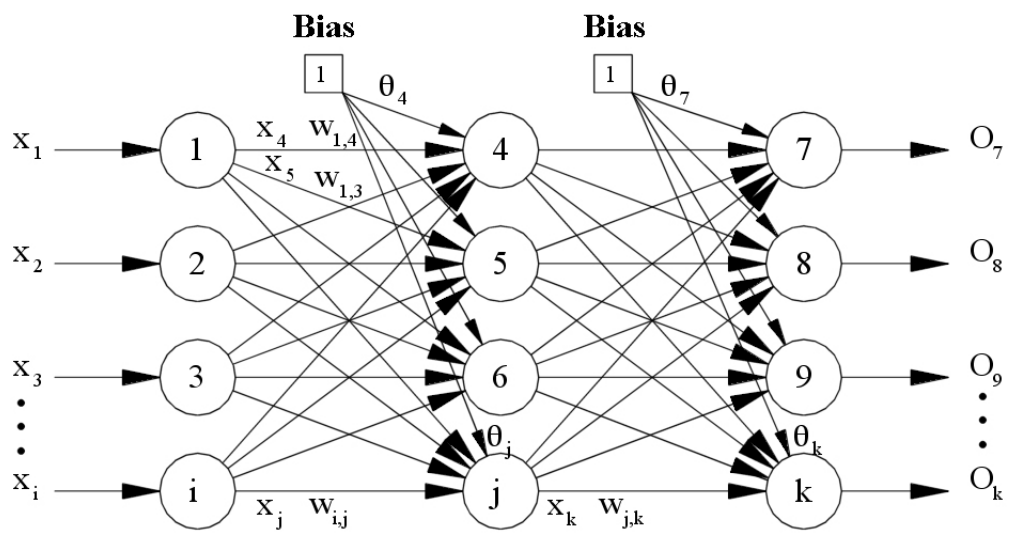


Bild 7.1: Beispielbild

# Bilderverzeichnis

- 2.1 Visualisierung grundlegender Prozesse während des Ablaufs eines überwachten Lernens . . . . . 6
- 2.2 Resultat eines auf unüberwachtem Lernen basierenden Clusterings . . . . . 7
- 2.3 Visualisierung grundlegender Prozesse während des Ablaufs eines unüberwachten Lernens . . . . . 7
- 5.1 Beschreibung der Grafik . . . . . 11
- 7.1 Beispielbild . . . . . 15

# Tabellen

7.1	Beispieltabelle . . . . .	14
-----	---------------------------	----

( Lorenz 2020)

## Literatur