

# Robustness & Graph (Convolutional) Neural Networks

---

Tim Bohne

25. Februar 2021

*Machine Learning Seminar 20/21*

## ① Motivation

## ② Graph Neural Networks (GNNs)

## ③ Robustheit

GNNs: Manipulation der Knoten-Attribute

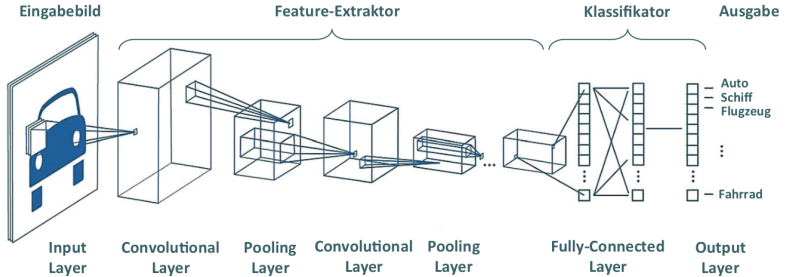
GNNs: Manipulation der Graph-Struktur

## ④ Fazit / Ausblick

## Motivationen für Graph Neural Networks:

- ① Convolutional Neural Networks (CNNs)
- ② Graph Embedding

# Convolutional Neural Networks (CNNs)



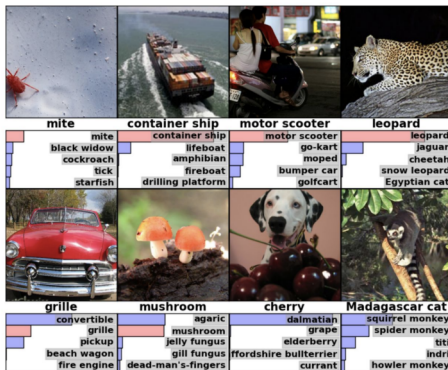
Convolutional Neural Network [9]

- **Feature-Learning**  
Sequenz aus Convolutional- und Pooling-Layern
- **Classification**

# Convolutional Neural Networks (CNNs)

IMAGENET

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.



ImageNet Challenge [6]

Hier sind CNNs sehr erfolgreich!

# Convolutional Neural Networks (CNNs)

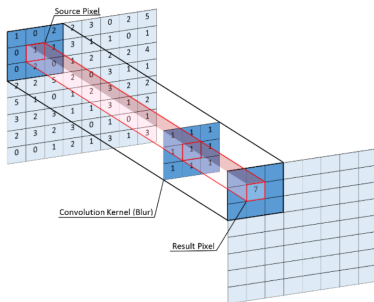
**Problem:** CNNs funktionieren lediglich mit Euklidischen Datenstrukturen (Bilder, Text), nicht mit Graphen!



107	108	114	100	100	100	120	101	172	145	100	100
100	140	142	74	74	43	81	17	110	210	100	104
100	140	50	14	34	4	14	34	40	100	140	141
204	100	4	134	107	111	122	204	100	14	140	141
104	64	107	401	407	100	100	100	107	47	71	101
170	100	107	203	214	200	200	204	14	14	104	141
64	170	200	100	210	211	100	107	14	14	104	141
100	47	140	84	14	100	100	11	41	42	14	141
100	104	141	100	100	107	170	140	101	14	101	141
203	174	180	102	100	207	140	170	100	45	104	141
100	214	174	140	204	107	84	104	74	14	204	141
140	144	107	100	107	110	111	140	100	100	140	141
100	214	174	140	100	141	141	141	141	141	141	141
140	141	100	141	141	141	141	141	141	141	141	141
140	141	100	141	141	141	141	141	141	141	141	141

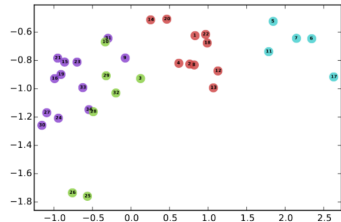
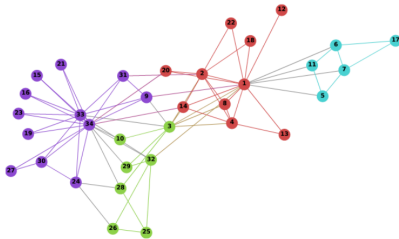
Pixelmatrix [8]

107	105	116	140	142	137	107	113	125	100	106
100	145	150	74	74	43	81	17	110	203	149
100	140	50	14	34	4	14	34	40	100	140
204	100	4	130	107	111	122	204	100	14	140
104	64	107	401	407	100	100	100	107	47	71
170	100	107	203	214	200	200	204	14	14	104
64	170	200	100	210	211	100	107	14	14	104
100	47	140	84	14	100	100	11	41	42	14
100	104	141	100	100	107	170	140	101	14	101
203	174	180	102	100	207	140	170	100	45	104
100	214	174	140	204	107	84	104	74	14	204
140	144	107	100	107	110	111	140	100	100	140
100	214	174	140	100	141	141	141	141	141	141
107	140	100	141	141	141	141	141	141	141	141
140	141	100	141	141	141	141	141	141	141	141
140	141	100	141	141	141	141	141	141	141	141



Convolution-Operation [4]

## Übersetze Graph-Struktur in niedrigdimensionalen Vektor



2-dim. Output für jeden Knoten [7]

## Ermöglichen Machine Learning auf Graphen, z.B.:

- Modellierung physikalischer Systeme
- Lernen molekularer Fingerabdrücke
- Analyse sozialer Netzwerke (Vorhersagen / Empfehlungen)
- Empfehlungssysteme



Freunde-Netzwerk einer Person [1]



## Typische Aufgaben für GNNs:

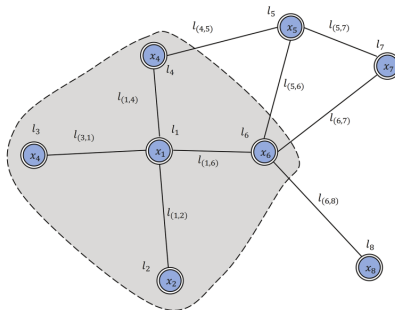
- Semi-Supervised Node Classification
- Graph Classification
- Link Prediction
- Clustering

# Graph Neural Networks (GNNs)

Graph  $G = (\mathbf{X}, \mathbf{A})$

$\mathbf{X}$  : Knoten - Personen in sozialem Netzwerk

$\mathbf{A}$  : Kanten - Beziehungen zwischen Personen



Graph basierend auf Scarselli et al. [5]

**Ziel:** Lerne „Node Embedding“  $\mathbf{h}_v$  für jeden Knoten  $v \in \mathbf{X}$  um einen Output  $\mathbf{o}_v$  zu generieren, z.B. das vorhergesagte Label

2 wichtige Funktionen ( $x$ : Input Feature,  $h$ : Hidden State):

① **Node Embedding:**  $h_v = f(x_v, x_{co[v]}, h_{ne[v]}, x_{ne[v]})$

② **Output Embedding:**  $o_v = g(h_v, x_v)$

Lerne Parameter von  $f$  und  $g$ :

**Loss Term:**  $\sum_{i=1}^P (t_i - o_i)$

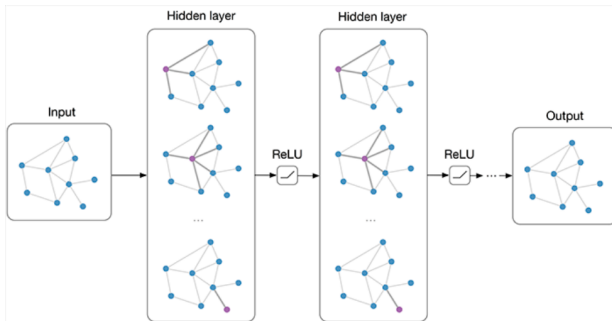
**Lernmethode:** Gradient Descent

- Iteratives Update:  $H^{t+1} = F(H^t, X)$
- Gradient wird basierend auf dem Loss berechnet
- Gewichte  $W$  werden basierend auf dem Gradienten adaptiert

# Graph Convolutional Networks (GCNs)

**Input:** Features  $X$ , Adjazenzmatrix  $A$

**Output:** Feature-Vektor für jeden Knoten



Multi-Layer GCN [3]

**Anschließend „Node Classification“ basierend auf den Output-Features**

Standard-Convolution kann nicht für Graphen definiert werden!

**Lösung:** Embeddings (spectral / spatial)

**Approximierte „Graph-Convolution“** (Kipf et al. [3]):

$$H^{(l+1)} = \sigma[\hat{D}^{-\frac{1}{2}} \hat{A} \hat{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}]$$

Sehr erfolgreich in praktischen Anwendungen eingesetzt!

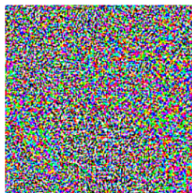
**Machine Learning Modelle sind i.A. anfällig für „Adversarial Attacks“**



**panda**

57.7% confidence

+ .007 ×



**noise**

=



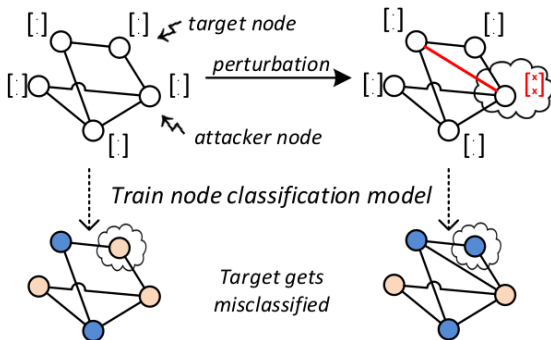
**gibbon**

99.3% confidence

Adaptiert aus Goodfellow et al. [2]

Dies gilt auch für GNNs, wobei sich zwei Arten von Angriffen unterscheiden lassen:

- **Manipulation der Knoten-Attribute**
- **Manipulation der Graph-Struktur**



## 3 Phasen in der Literatur:

- ① GNNs anfällig für „Adversarial Attacks“
- ② Verteidigungsmechanismen für konkrete Szenarien
- ③ Beweisbare Garantien für die Robustheit bestimmter Modelle



# Semi-Supervised Node Classification

Graph  $\mathbf{G} = (\mathbf{A}, \mathbf{X})$     Zielknoten  $\mathbf{t}$     Trainierbare Parameter  $\theta$

$\mathbf{A}$ : Adjazenzmatrix,  $\mathbf{X}$ : Knoten-Features

$\mathbf{V}$ : Knotenmenge

$\mathbf{V}_L \subseteq \mathbf{V}$ : Teilmenge der gelabelten Knoten

$\mathcal{T}(\mathbf{A})$ : „Message-Passing-Matrix“ - wie Aktivierungen durch das GNN propagiert werden  $\rightarrow$  Transformation der Adjazenzmatrix

**Cross-Entropy-Loss Minimierung:**

Lerne  $\theta$  unter Verwendung der  $\mathbf{V}_L$

**Ziel:** Label der verbleibenden Knoten ohne Label vorhersagen

## **Zertifizierte Robustheit gegenüber Manipulationen der Knoten-Attribute** - Zügner et al. [11]

Wie lässt sich sicherstellen, dass kleine Änderungen an den Attributen keine dramatischen Auswirkungen auf den Output haben?

**Ziel:** Zertifikat für Knoten  $t$  bedeutet, dass die Vorhersage für  $t$  sich nach zulässigen Manipulationen nicht ändert

→ Definiertes Angriffsmodell

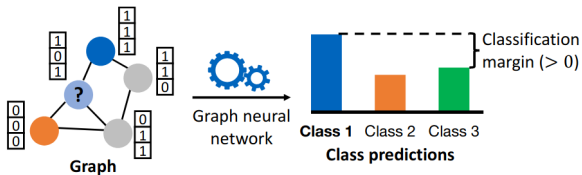
**Idee:** „Worst-Case-Margin“  $m^t$  für Knoten  $t$  zwischen den Klassen  $y$  und  $y^*$  bzgl. zulässiger Manipulationen:

$$m^t(y^*, y) := \min_{\tilde{X}} f_{\theta}^t(\tilde{X}, A)_{y^*} - f_{\theta}^t(\tilde{X}, A)_y$$

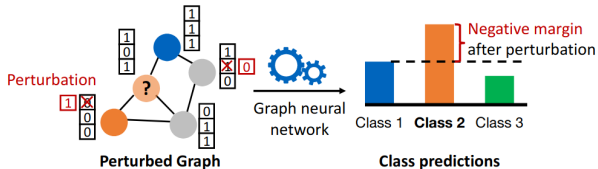
**s.t.**  $\tilde{X} \in \text{zul. Manipulationen}$

$m^t > 0 \forall y \neq y^* \rightarrow$  Modell robust bzgl. Knoten  $t$

# High-Level Idee



Vor der Manipulation [11]



Nach der Manipulation [11]

## Optimierungsproblem nicht effizient lösbar:

Diskrete Daten + nicht-konvexe Aktivierungsfunktion

## Lösung:

- Konvexe ReLU Relaxation: Transformiert das GNN in ein effizient lösbares LP
- Kontinuierliche Relaxation der ganzzahligen Knotenattribute

→ Effizient berechenbare Schranken für den „Worst-Case-Margin“

→ Ggf. „false negatives“, jedoch keine „false positives“

**Ziel:** Erreiche beweisbare Robustheit durch Training

→ Optimierte bezüglich Robustheit

## Robust Hinge Loss

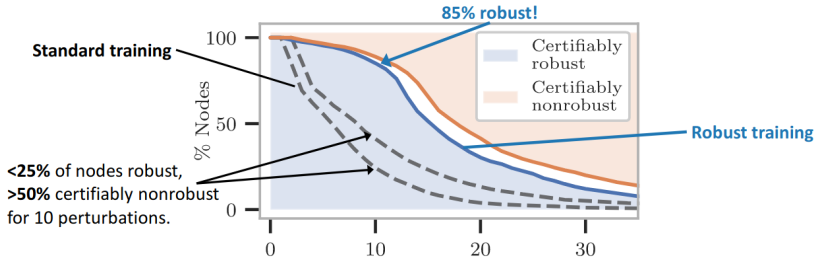
$$\min_{\theta} \sum_{v \in V_L} \mathcal{L}(p_v, y_v) + \sum_{v \in V_L} \hat{\mathcal{L}}_{M_L}(-m_v^*, y_v) + \sum_{v \in V \setminus V_L} \hat{\mathcal{L}}_{M_U}(-m_v^*, \hat{y}_v)$$

- Relaxierte Version um Robustheit sicherzustellen
- Exakte Version für die Klassifizierung

→ **Mit Standardsoftware robuste GNNs trainieren**

→ **Qualität der Klassifizierungsergebnisse nicht schlechter**

# Training mit dem Ziel der Robustheit

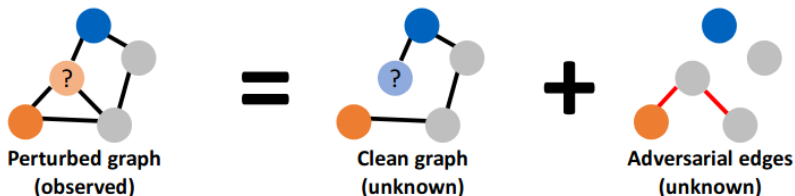


„Robustheits-Training“ vs. Standard-Training [11]

- Nur wenige Knoten nicht zertifizierbar
- Standard-Training führt zu GNNs, die nur robust bzgl. sehr weniger Manipulationen sind

## Zertifizierte Robustheit gegenüber Manipulationen der Graph-Struktur - Zügner et al. [12]

Angreifer können neue Kanten in den Graphen einfügen, z.B. Likes in Social Media:



High-Level Idee [12]



**Einfügen neuer Kanten kann die Vorhersagen des Modells ändern!**

**Ziel:** Ein Zertifikat bedeutet, dass sich die Vorhersage für einen Knoten nicht ändert, wenn ein Angreifer Kanten hinzufügt

**Idee:** Prüfe, ob es einen Graphen gibt, der durch Entfernen von Kanten erreichbar ist und die Vorhersage ändert

**Angriffsmodell:** Praktisch unmöglich, alle Graphen zu enumerieren → Limitiere Anzahl an „Adversarial Edges“

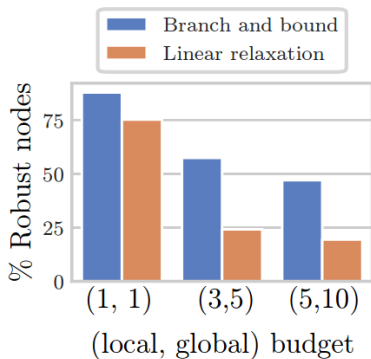
- Bereits trainiertes GNN mit Parametern  $\theta$
- Input-Graph ( $\mathbf{A}$ ) möglicherweise manipuliert
- $\mathbf{A}$  ist von der „sauberen“ Variante  $\mathbf{A}^*$  durch eine Menge zulässiger Manipulationen erreichbar

**Worst-Case-Margin:**

$$m^t(y^*, y) := \min_{\tilde{\mathbf{A}}} f_{\theta}^t(X, \tilde{\mathbf{A}})_{y^*} - f_{\theta}^t(X, \tilde{\mathbf{A}})_y$$
$$\text{s.t. } \tilde{\mathbf{A}} \in \text{zul. Manipulationen}$$

$m^t > 0 \forall y \neq y^* \rightarrow$  Modell robust bzgl. Knoten  $t$

- Effiziente Zertifizierung der Robustheit gegenüber Manipulationen der Graph-Struktur
- Bereits wenige Manipulationen führen zu einem Label-Wechsel eines beträchtlichen Anteils der Knoten



- GNNs sind sehr **erfolgreich** in einer Vielzahl von praktischen Anwendungen einsetzbar
- GNNs sind **nicht robust** - anfällig für Manipulationen
  - der Knoten-Attribute
  - der Graph-Struktur
- Um GNNs in (sicherheitskritischen) praktischen Anwendungen nutzen zu können, ist ein bestimmtes Maß an **Robustheit** **nötig**
- **Beweisbare Garantien** für die (Nicht-)Robustheit von GNNs sind möglich und nötig - erste Ansätze gesehen

- Spielraum für **Generalisierungen** (z.B. Angriffsmodelle)
- Trainingsmethoden, die **Robustheit explizit als Optimierungsziel** enthalten
- **Robustheitsgarantien** für Attribut- und Struktur-Manipulationen
- **Grundsätzliches Verständnis** - was macht Manipulationen schädlich?

## Ideen / Ansätze

- Graphen, die reale Probleme repräsentieren, besitzen bestimmte strukturelle Gemeinsamkeiten, die manipulierte Graphen verletzen
- Statistisch signifikante strukturelle Attribute in Manipulationen entdecken
- Wahrscheinlichkeit dafür, dass bestimmte Manipulationen schädlich sind

→ **Allgemeinere Angriffsmodelle**

→ **Manipulationen erkennen + verhindern**

→ **Robustheit**



Yihui Fan.

**Creating and Analysing Facebook Friend Network Graphs Using Python.**

[https:](https://www.databentobox.com/2019/07/28/facebook-friend-graph/)

[//www.databentobox.com/2019/07/28/facebook-friend-graph/](https://www.databentobox.com/2019/07/28/facebook-friend-graph/).

Abgerufen am: 15.02.2021.



Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy.

**Explaining and harnessing adversarial examples.**

In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, San Diego, CA, USA, May 7-9, Conference Track Proceedings*, 2015.



Thomas N. Kipf and Max Welling.

**Semi-supervised classification with graph convolutional networks.**

In *5th International Conference on Learning Representations, Toulon, France, April 24-26, Conference Track Proceedings*. OpenReview.net, 2017.



Ahmad Kiswani.

**Convolution and Normalized Cross Correlation on Kepler Architecture.**

<https://simpl.eelabs.technion.ac.il/projects/convolution-and-normalized-cross-correlation-on-kepler-architecture>

Abgerufen am: 15.02.2021.



Zhiyuan Liu and Jie Zhou.

**Introduction to graph neural networks.**

*Synthesis Lectures on Artificial Intelligence and Machine Learning*,  
14:1–127, 03 2020.



X. Giro o Nieto.

**Image classification on Imagenet (D1L4 2017 UPC Deep Learning for Computer Vision).**

<https://www.slideshare.net/xavigiro/image-classification-on-imagenet-d1l4-2017-upc-deep-learning-for-c>

Abgerufen am: 15.02.2021.





Bryan Perozzi, Rami Al-Rfou, and Steven Skiena.

**Deepwalk.**

*Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, Aug 2014.



Thomas Smits and Melvin Wevers.

**The visual digital turn: Using neural networks to study historical images.**

*Digital Scholarship in the Humanities*, 35, 01 2018.



Patrick Zschech, Christoph Sager, Philipp Siebers, and Maik Pertermann.

**Mit computer vision zur automatisierten qualitätssicherung in der industriellen fertigung: Eine fallstudie zur klassifizierung von fehlern in solarzellen mittels elektrolumineszenz-bildern.**

*HMD Praxis der Wirtschaftsinformatik*, 07 2020.



Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann.

**Adversarial attacks on neural networks for graph data.**

In Yike Guo and Faisal Farooq, editors, *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, London, UK, August 19-23*, pages 2847–2856. ACM, 2018.



Daniel Zügner and Stephan Günnemann.

**Certifiable robustness and robust training for graph convolutional networks.**

In Ankur Teredesai, Vipin Kumar, Ying Li, Rómer Rosales, Evimaria Terzi, and George Karypis, editors, *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, Anchorage, AK, USA, August 4-8*, pages 246–256. ACM, 2019.



Daniel Zügner and Stephan Günnemann.

**Certifiable robustness of graph convolutional networks under structure perturbations.**

In Rajesh Gupta, Yan Liu, Jiliang Tang, and B. Aditya Prakash, editors, *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27*, pages 1656–1665. ACM, 2020.