

---

# EVALUATION VON LERNVERFAHREN ZUR BEWERTUNG DER GEFAHRENLAGE IM INTERNET

Diplomarbeit an der Universität Bonn, Institut für Informatik IV

Simon Hunke

---

---

# Agenda

---

1. Einführung
2. Methodik & untersuchte Lernverfahren
3. Konkretisierung des Lernproblems
4. Resultate & Fazit
5. Ausblick & Empfehlungen

# Automatisierte Gefahrenbewertung der Gesamtlage

- **Motivation:** Unterstützung von Experten
  - Angriffsaufkommen im Internet übersteigt menschliche Kapazität
  - Erfahrung der Uni Bonn: Mehrere tausend Angriffe pro Tag und Host
  
- Strategien: regelbasiert oder Lernverfahren
  - Maschinelles Lernen ist interessant wegen ...
    - ... komplexer Problemdomäne
    - ... hoher Dynamik der Technologien im Internet

## **Ziel der Diplomarbeit:**

Quantitative Evaluation von Lernverfahren

# Beispiele zu verwandten Arbeiten

Anomalieerkennung mit Lernverfahren:

- Intrusion Detection System auf Basis neuronaler Netze
- Spam-Filter mit Bayes Netzen
- Malware-Identifikation mit Support Vector Machines
- ➡ Abgrenzung zur Diplomarbeit: 0-1-Entscheidungen statt Abstufung

Internet-Lagebewertung:

- Internet Storm Center
- Symantec ThreatCon
- ➡ Abgrenzung zur Diplomarbeit: Manuelle Kategorisierung der Lage

# Leistungsbewertung der Lernverfahren

**Kontext:** Überwachtes Lernen

- Generalisierungsvermögen einer Hypothese entscheidend

Untersuchte Lernverfahren:

- Neuronale Netze & Bayes Netze

Untersuchte Metriken zur Leistungsbewertung der Lernverfahren:

- *Präzision* (Prozentsatz korrekt klassifizierter Beispiele)
- *Lerndauer* (Zeit pro Lernvorgang)
- *Steigerungspotential* (in Bezug auf Präzision)

Generelle Methodik zur Präzisionsschätzung:

- Aufteilen der Beispiele in eine *Trainingsmenge* und eine *Testmenge*

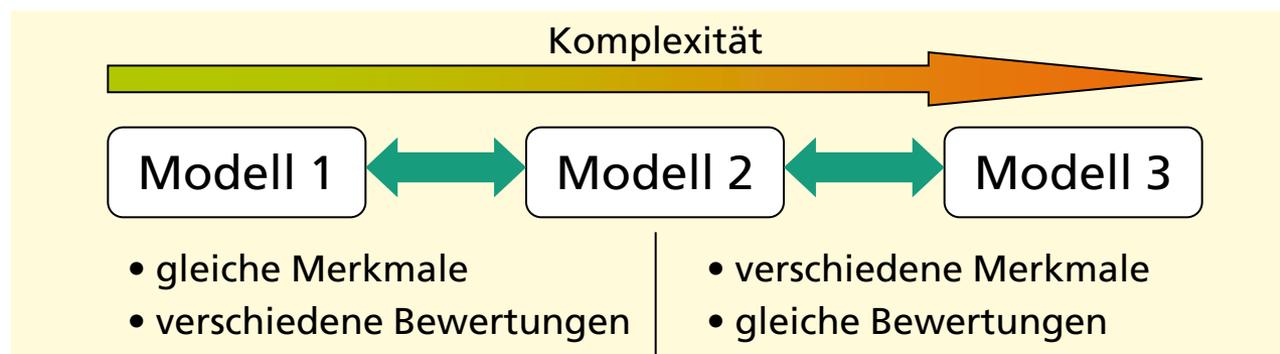
# Internet-Lagebewertung als Lernproblem

Betrachteter Modellrahmen: (Uni Bonn)

## ■ Modellaufbau:

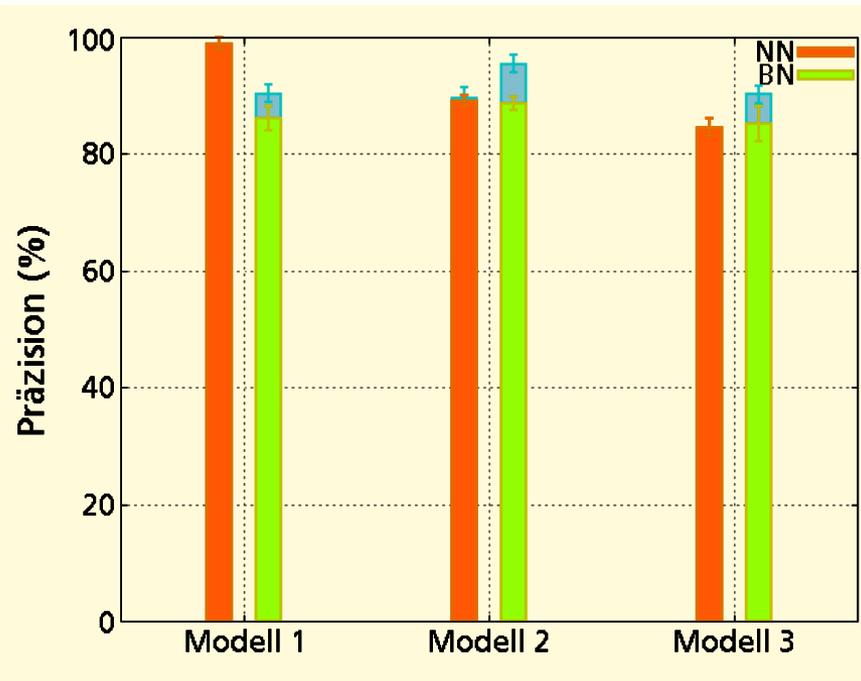
- Lagebild (endliche Merkmalsmenge zur Lagebeschreibung)
- Bewertungsfunktion (Simulation von Experteneinschätzungen)

## ■ Fokus: Einfluss von Komplexität auf Performanz von Lernverfahren

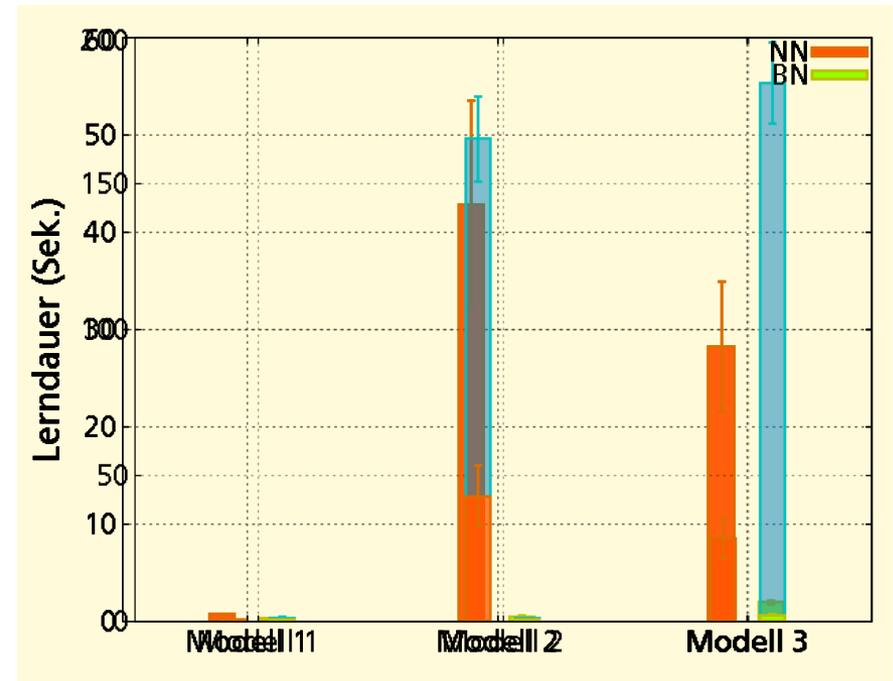


# Präzision & Lerndauer

Präzision:

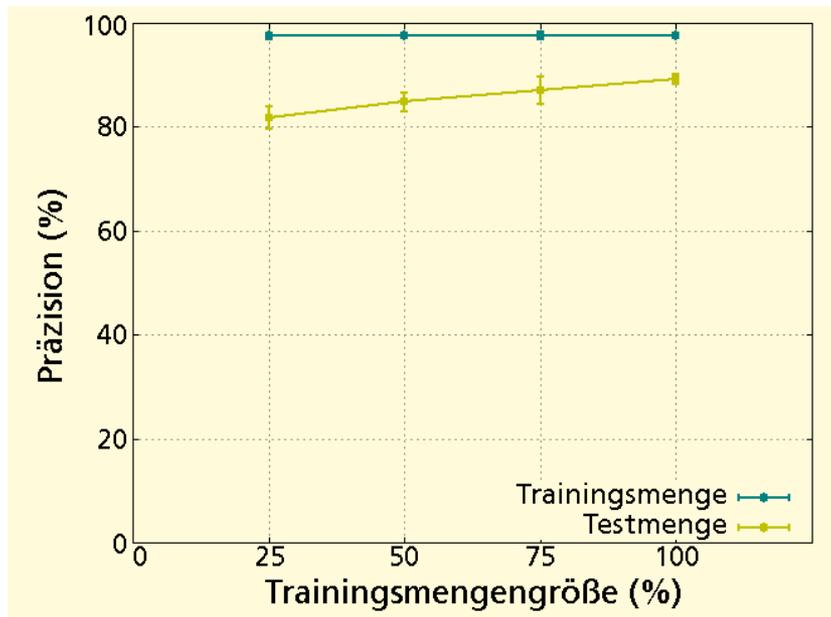


Lerndauer:



# Einfluss der Trainingsmengengröße

## Neuronale Netze:



**Bayes Netze:** Erwartungsgemäß keine Steigerung

## „Curse of Dimensionality“

Bei hochdimensionalen Eingaberaum können „Löcher“ in den Trainingsdaten entstehen.



wichtige Infos für hohe Generalisierung fehlen!

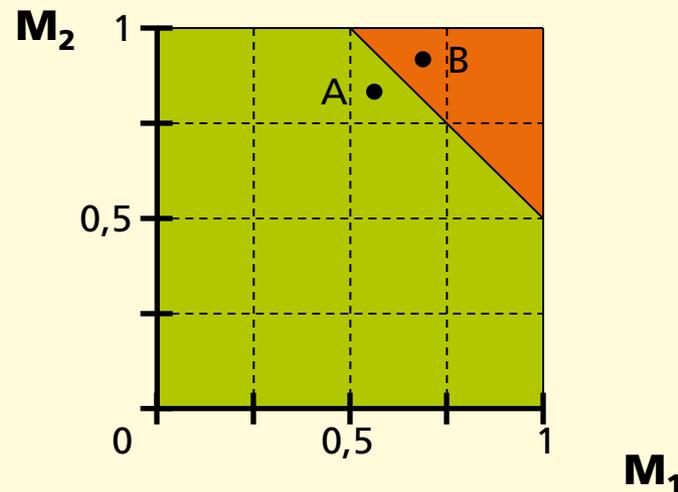
**Aber:** Auch Komplexität hat Einfluss!

# Bayes Netze

## Lernfaktor Diskretisierung

Lernalgorithmen für Bayes Netze unterstützen nur diskrete Variablen

➔ durch Diskretisierung entsteht *Fehler* !

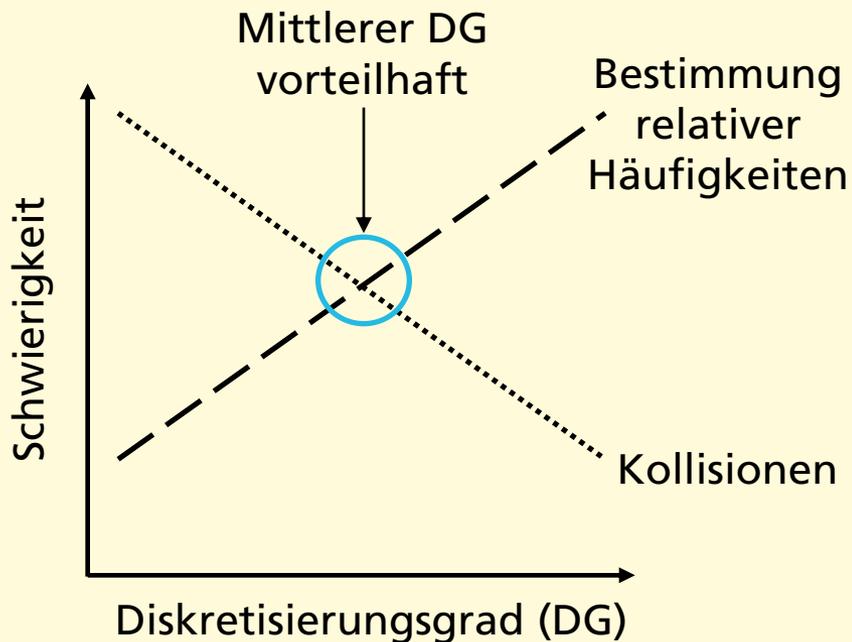


Jedes Bayes Netz kann entweder nur A oder B korrekt erkennen !

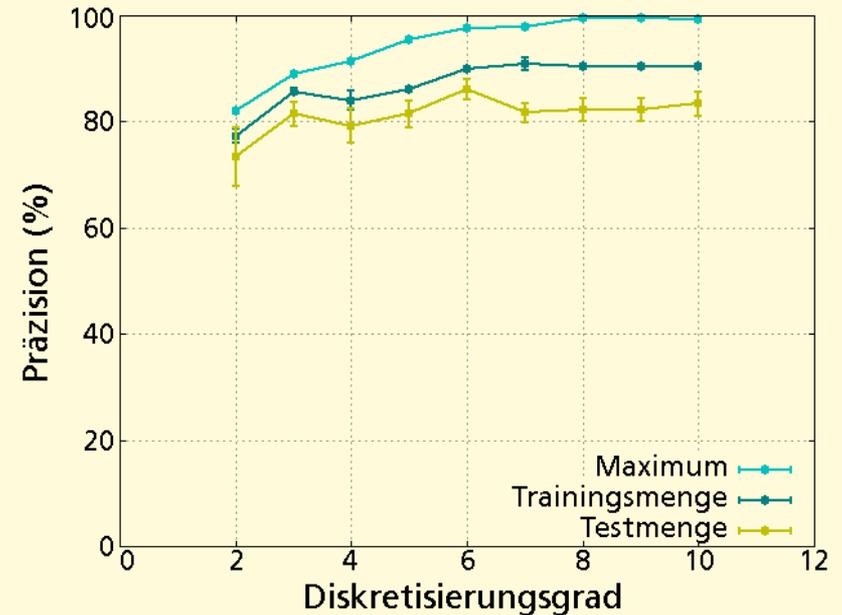
# Bayes Netze

## Lernfaktor Diskretisierung (2)

### Theorie:



### Auswirkungen in der Praxis:



# Fazit

- Beide Lernverfahren auch bei hoher Modellkomplexität gute Leistung
  
- Neuronale Netze
  - Benötigen bei steigender Modellkomplexität mehr Zeit und Daten
  - Präzision korreliert negativ mit Modellkomplexität (bei gleichem Datenvolumen)
  - Nahezu unabhängig von der Art der Modellierung (kont. oder disk.)
  
- Bayes Netze
  - Eignen sich am besten für diskrete Modelle
  - Greedy-Heuristiken erlauben schnelle Lernzeiten

# Ausblick & Empfehlungen

Weiterführende Untersuchungen:

- Evaluation weiterer Lernverfahren
  - z.B. Entscheidungsbäume, Clustering-Verfahren, ...
- Konzipierung praxisrelevanter Modelle
  - Auch reale Experteneinschätzungen (Inkonsistenzen?)

Empfehlungen für die Praxis:

- Analyse der Trainingsdaten
- Aufbereitung der Trainingsdaten
  - z.B. möglichst verlustfreie Diskretisierung, Integration von verfügbarem Wissen

**Danke für die Aufmerksamkeit !**

**Fragen?**