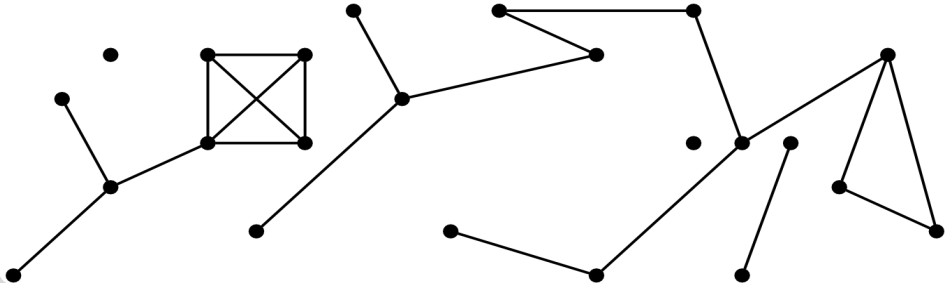
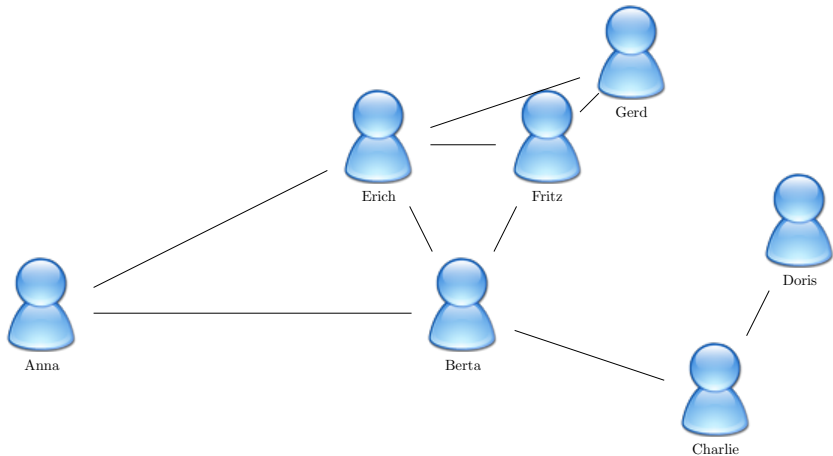


# On Node Classification in Dynamic Content-based Networks

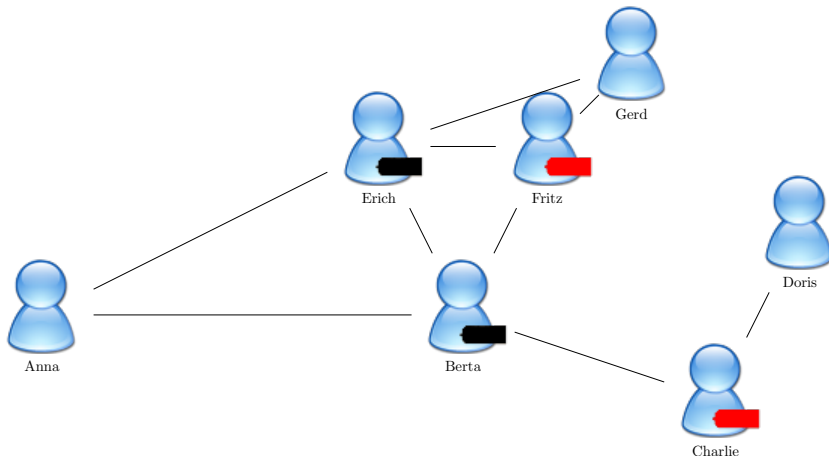
Martin Thoma | 28. Februar 2014

INSTITUT FÜR PROGRAMMSTRUKTUREN UND DATENORGANISATION

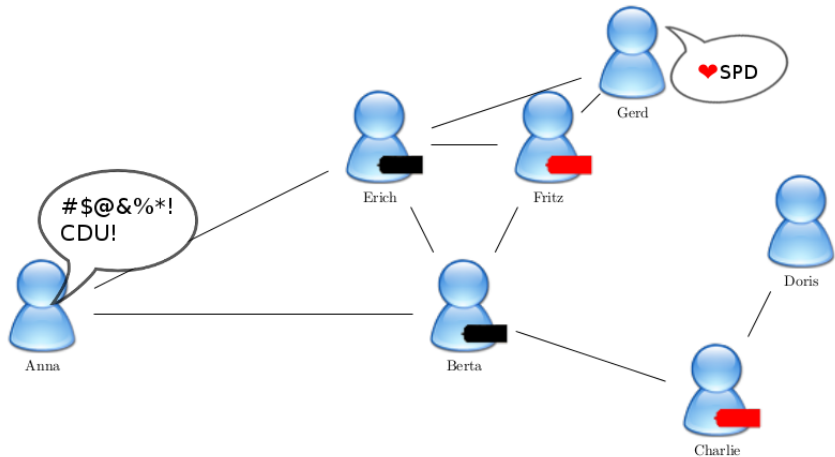




# Partially labeled network



# Partially labeled network with content



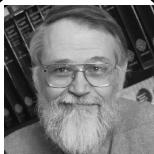
# Beispiel 2: Literaturdatenbanken

The Development  
of the C Language  
Interprocess  
Communication in  
the Ninth Edition  
Unix System



Computer Science

The C Programming  
Language  
digital restoration  
and typesetter

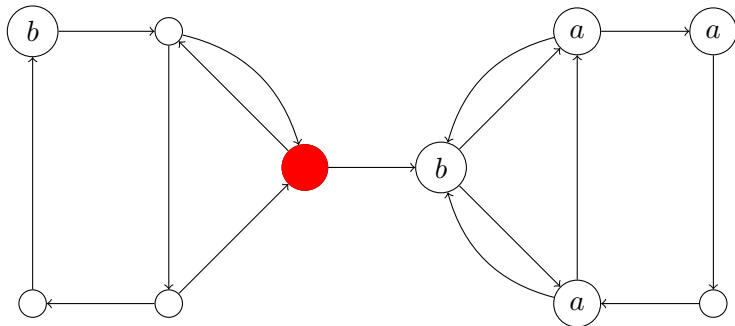


Computer Science

The Identity  
Thesis for  
Language and  
Music

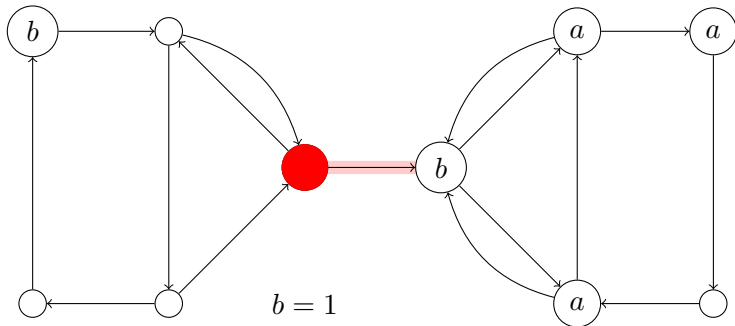


Linguistics



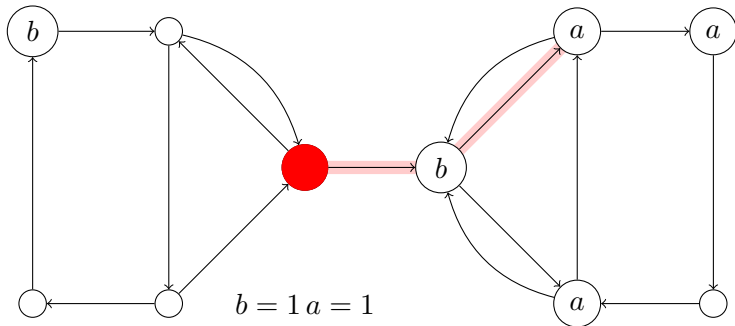
Klassifizieren des roten Knotens:

- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren



Klassifizieren des roten Knotens:

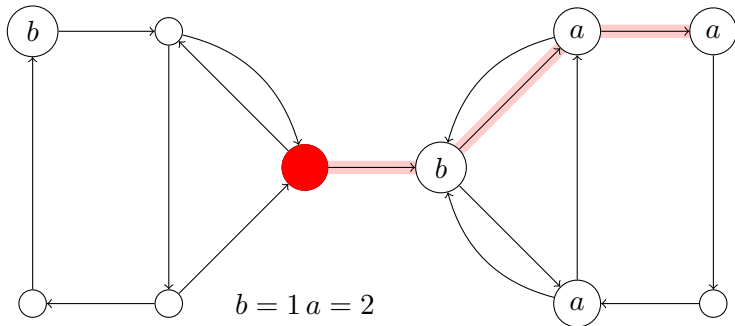
- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren



Klassifizieren des roten Knotens:

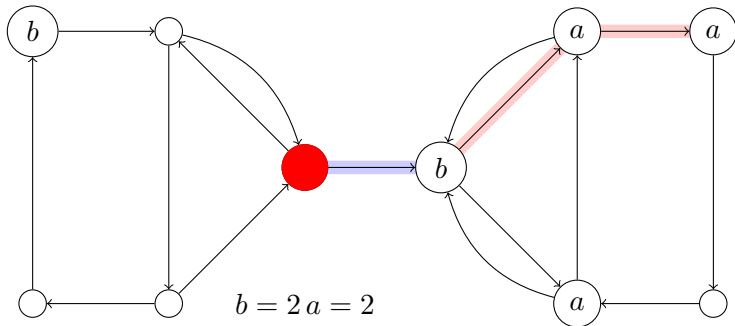
- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren





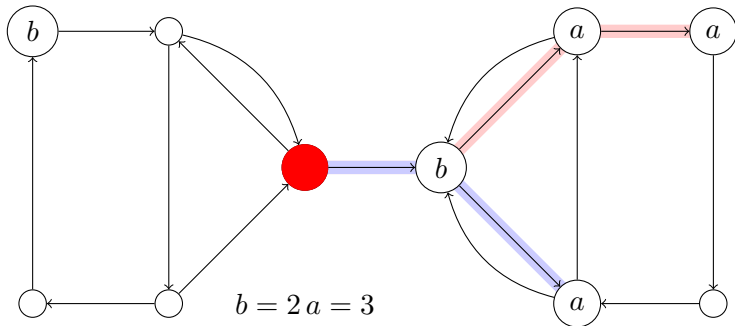
Klassifizieren des roten Knotens:

- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren



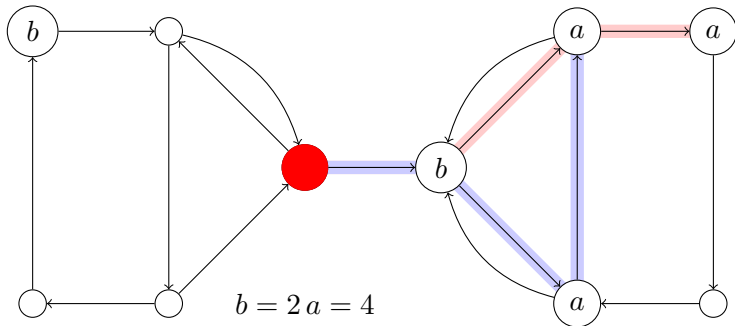
Klassifizieren des roten Knotens:

- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren



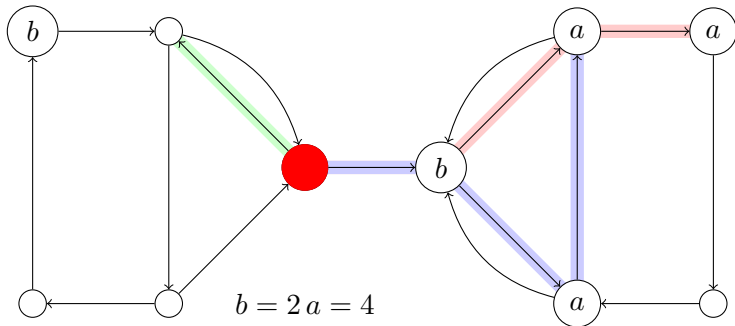
Klassifizieren des roten Knotens:

- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren



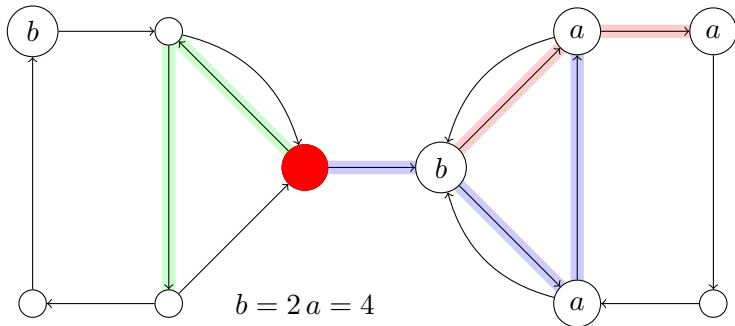
Klassifizieren des roten Knotens:

- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren

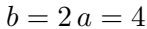


Klassifizieren des roten Knotens:

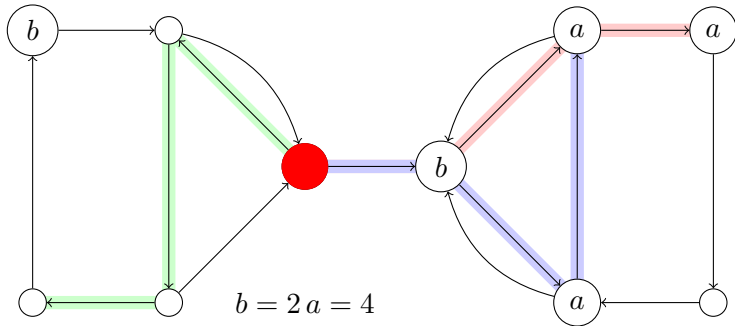
- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren



- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren



- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren



Klassifizieren des roten Knotens:

- Zählen von Knotenbeschriftungen in Random Walks
- 4 Random Walks, beginnend bei Rot
- 3 Sprünge pro Random Walk
- $4 \cdot a, 2 \cdot b \Rightarrow$  Rot mit  $a$  klassifizieren



- Neben Struktur können Texte genutzt werden
- Einschränkung: Effizienz!
- Idee: Graph erweitern
  - Texte als Wortmengen
  - Strukturknoten verweisen auf Wortknoten
  - vice versa

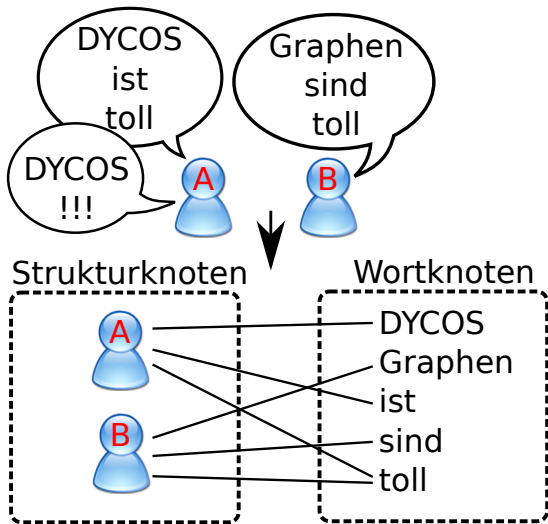
- Neben Struktur können Texte genutzt werden
- Einschränkung: Effizienz!
- Idee: Graph erweitern
  - Texte als Wortmengen
  - Strukturknoten verweisen auf Wortknoten
  - vice versa

- Neben Struktur können Texte genutzt werden
- Einschränkung: Effizienz!
- Idee: Graph erweitern
  - Texte als Wortmengen
  - Strukturknoten verweisen auf Wortknoten
  - vice versa

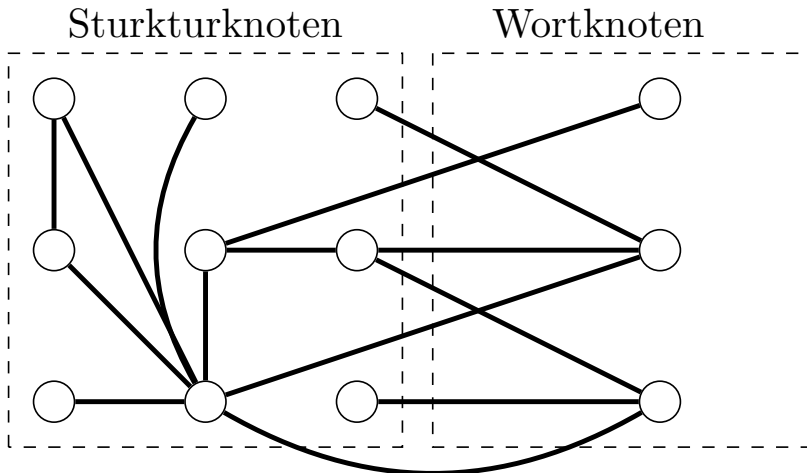
- Neben Struktur können Texte genutzt werden
- Einschränkung: Effizienz!
- Idee: Graph erweitern
  - Texte als Wortmengen
  - Strukturknoten verweisen auf Wortknoten
  - vice versa

- Neben Struktur können Texte genutzt werden
- Einschränkung: Effizienz!
- Idee: Graph erweitern
  - Texte als Wortmengen
  - Strukturknoten verweisen auf Wortknoten
  - vice versa

- Neben Struktur können Texte genutzt werden
- Einschränkung: Effizienz!
- Idee: Graph erweitern
  - Texte als Wortmengen
  - Strukturknoten verweisen auf Wortknoten
  - vice versa



# Erweiterter, semi-bipartiter Graph





- Füllwörter: und, oder, im, in, ...

⇒ Beschränkung des Vokabulars sinnvoll

Idee:

- Zufällige Beispielmenge von Texten für Vokabularbildung betrachten
- Gini-Koeffizient nutzen

- Füllwörter: und, oder, im, in, ...

⇒ Beschränkung des Vokabulars sinnvoll

Idee:

- Zufällige Beispielmenge von Texten für Vokabularbildung betrachten
- Gini-Koeffizient nutzen

- Füllwörter: und, oder, im, in, ...

⇒ Beschränkung des Vokabulars sinnvoll

Idee:

- Zufällige Beispielmenge von Texten für Vokabularbildung betrachten
- Gini-Koeffizient nutzen

- Füllwörter: und, oder, im, in, ...

⇒ Beschränkung des Vokabulars sinnvoll

## Idee:

- Zufällige Beispielmenge von Texten für Vokabularbildung betrachten
- Gini-Koeffizient nutzen

- Füllwörter: und, oder, im, in, ...

⇒ Beschränkung des Vokabulars sinnvoll

## Idee:

- Zufällige Beispielmenge von Texten für Vokabularbildung betrachten
- Gini-Koeffizient nutzen

- Füllwörter: und, oder, im, in, ...

⇒ Beschränkung des Vokabulars sinnvoll

## Idee:

- Zufällige Beispielmenge von Texten für Vokabularbildung betrachten
- Gini-Koeffizient nutzen

- statistisches Maß für Ungleichverteilung

- $g = \sum_i p_i^2$  mit  $p_i$  als relative Häufigkeit

- $g \in (0, 1]$

- $g$  nahe bei 1  $\Rightarrow$  Wort ist stark ungleich verteilt

$\Rightarrow$  Nehme Top- $m$  Wörter mit höchstem Gini-Koeffizient

- statistisches Maß für Ungleichverteilung
  - $g = \sum_i p_i^2$  mit  $p_i$  als relative Häufigkeit
  - $g \in (0, 1]$
  - $g$  nahe bei 1  $\Rightarrow$  Wort ist stark ungleich verteilt
- $\Rightarrow$  Nehme Top- $m$  Wörter mit höchstem Gini-Koeffizient



- statistisches Maß für Ungleichverteilung

- $g = \sum_i p_i^2$  mit  $p_i$  als relative Häufigkeit

- $g \in (0, 1]$

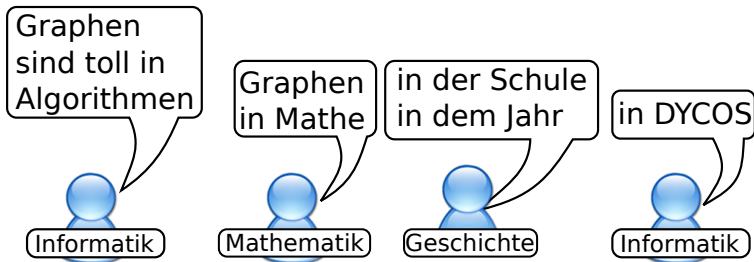
- $g$  nahe bei 1  $\Rightarrow$  Wort ist stark ungleich verteilt

$\Rightarrow$  Nehme Top- $m$  Wörter mit höchstem Gini-Koeffizient

- statistisches Maß für Ungleichverteilung
- $g = \sum_i p_i^2$  mit  $p_i$  als relative Häufigkeit
- $g \in (0, 1]$
- $g$  nahe bei 1  $\Rightarrow$  Wort ist stark ungleich verteilt

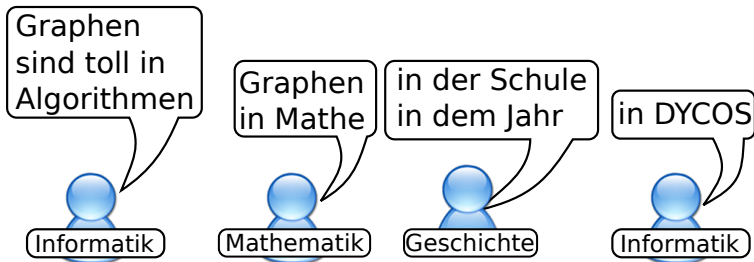
$\Rightarrow$  Nehme Top- $m$  Wörter mit höchstem Gini-Koeffizient

- statistisches Maß für Ungleichverteilung
  - $g = \sum_i p_i^2$  mit  $p_i$  als relative Häufigkeit
  - $g \in (0, 1]$
  - $g$  nahe bei 1  $\Rightarrow$  Wort ist stark ungleich verteilt
- $\Rightarrow$  Nehme Top- $m$  Wörter mit höchstem Gini-Koeffizient



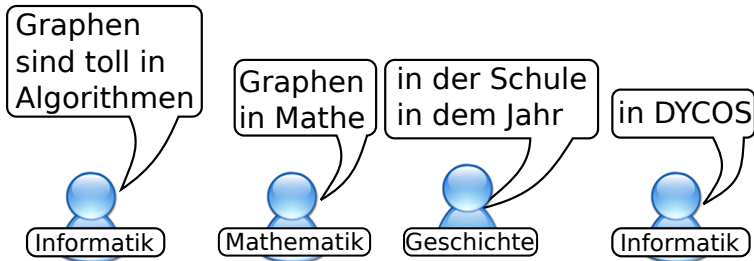
Beispiel: „in“

- Vorkommen insgesamt:  $5 \times$
- Vorkommen in „Informatik“  $2 \times \Rightarrow p_1 = \frac{2}{5}$
- Vorkommen in „Mathematik“  $1 \times \Rightarrow p_2 = \frac{1}{5}$
- Vorkommen in „Geschichte“  $2 \times \Rightarrow p_3 = \frac{2}{5}$
- Gini-Koeffizient:  $\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{2}{5}\right)^2 = \frac{9}{25}$



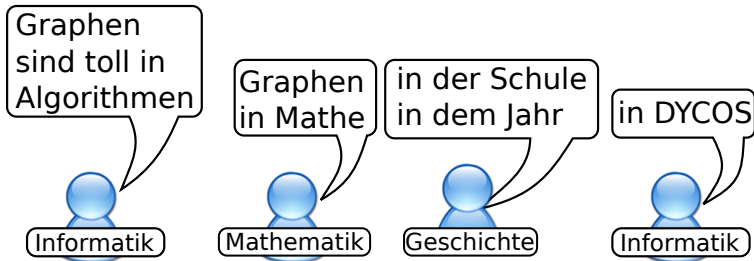
Beispiel: „in“

- Vorkommen insgesamt:  $5 \times$
- Vorkommen in „Informatik“  $2 \times \Rightarrow p_1 = \frac{2}{5}$
- Vorkommen in „Mathematik“  $1 \times \Rightarrow p_2 = \frac{1}{5}$
- Vorkommen in „Geschichte“  $2 \times \Rightarrow p_3 = \frac{2}{5}$
- Gini-Koeffizient:  $\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{2}{5}\right)^2 = \frac{9}{25}$



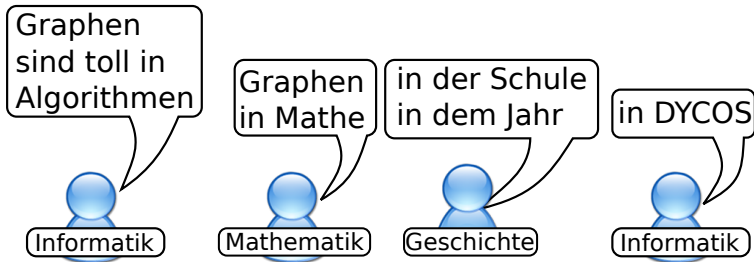
Beispiel: „in“

- Vorkommen insgesamt:  $5 \times$
- Vorkommen in „Informatik“  $2 \times \Rightarrow p_1 = \frac{2}{5}$
- Vorkommen in „Mathematik“  $1 \times \Rightarrow p_2 = \frac{1}{5}$
- Vorkommen in „Geschichte“  $2 \times \Rightarrow p_3 = \frac{2}{5}$
- Gini-Koeffizient:  $\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{2}{5}\right)^2 = \frac{9}{25}$



Beispiel: „in“

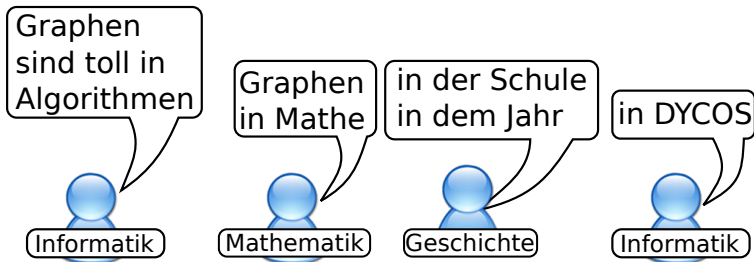
- Vorkommen insgesamt:  $5 \times$
- Vorkommen in „Informatik“  $2 \times \Rightarrow p_1 = \frac{2}{5}$
- Vorkommen in „Mathematik“  $1 \times \Rightarrow p_2 = \frac{1}{5}$
- Vorkommen in „Geschichte“  $2 \times \Rightarrow p_3 = \frac{2}{5}$
- Gini-Koeffizient:  $\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{2}{5}\right)^2 = \frac{9}{25}$



Beispiel: „in“

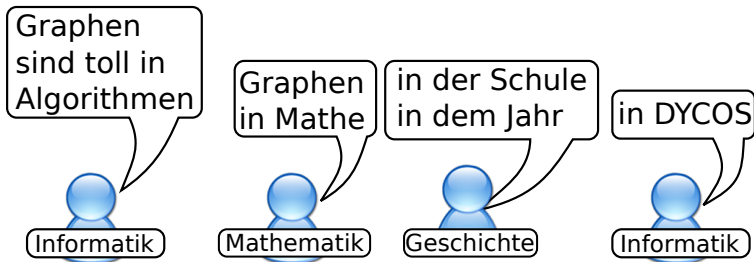
- Vorkommen insgesamt:  $5 \times$
- Vorkommen in „Informatik“  $2 \times \Rightarrow p_1 = \frac{2}{5}$
- Vorkommen in „Mathematik“  $1 \times \Rightarrow p_2 = \frac{1}{5}$
- Vorkommen in „Geschichte“  $2 \times \Rightarrow p_3 = \frac{2}{5}$
- Gini-Koeffizient:  $\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{2}{5}\right)^2 = \frac{9}{25}$





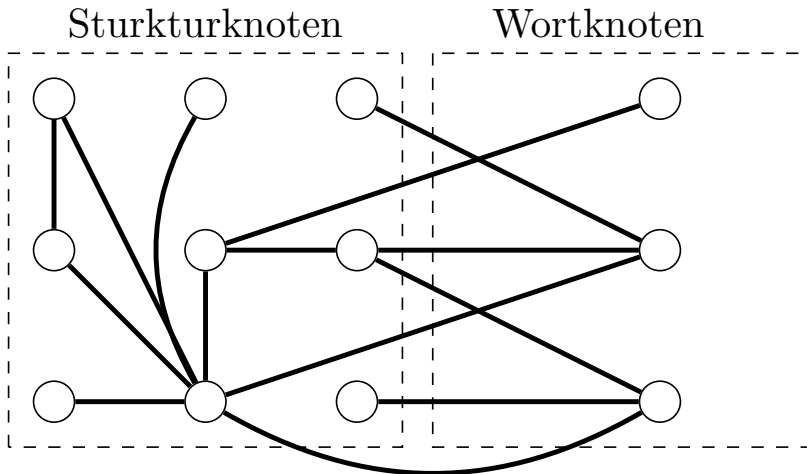
Beispiel: „in“

- Vorkommen insgesamt:  $5 \times$
- Vorkommen in „Informatik“  $2 \times \Rightarrow p_1 = \frac{2}{5}$
- Vorkommen in „Mathematik“  $1 \times \Rightarrow p_2 = \frac{1}{5}$
- Vorkommen in „Geschichte“  $2 \times \Rightarrow p_3 = \frac{2}{5}$
- Gini-Koeffizient:  $\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{2}{5}\right)^2 = \frac{9}{25}$



Beispiel: „in“

- Vorkommen insgesamt:  $5 \times$
- Vorkommen in „Informatik“  $2 \times \Rightarrow p_1 = \frac{2}{5}$
- Vorkommen in „Mathematik“  $1 \times \Rightarrow p_2 = \frac{1}{5}$
- Vorkommen in „Geschichte“  $2 \times \Rightarrow p_3 = \frac{2}{5}$
- Gini-Koeffizient:  $\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{2}{5}\right)^2 = \frac{9}{25}$



- **Struktursprung:** von Strukturknoten  $v$  zu Strukturknoten  $v'$
- **Inhaltlicher Zweifachsprung:** von Strukturknoten  $v$  über Wortknoten zu Strukturknoten  $v'$ 
  - Finde alle Knoten  $v'$ , die über Wortknoten erreichbar sind (Pfadlänge 2)
  - Nehme Top- $q$ -Knoten (Anzahl der Pfade)
  - Wähle zufällig einen davon

- **Struktursprung:** von Strukturknoten  $v$  zu Strukturknoten  $v'$
- **Inhaltlicher Zweifachsprung:** von Strukturknoten  $v$  über Wortknoten zu Strukturknoten  $v'$ 
  - Finde alle Knoten  $v'$ , die über Wortknoten erreichbar sind (Pfadlänge 2)
  - Nehme Top- $q$ -Knoten (Anzahl der Pfade)
  - Wähle zufällig einen davon

- **Struktursprung:** von Strukturknoten  $v$  zu Strukturknoten  $v'$
- **Inhaltlicher Zweifachsprung:** von Strukturknoten  $v$  über Wortknoten zu Strukturknoten  $v'$ 
  - Finde alle Knoten  $v'$ , die über Wortknoten erreichbar sind (Pfadlänge 2)
  - Nehme Top- $q$ -Knoten (Anzahl der Pfade)
  - Wähle zufällig einen davon

- **Struktursprung:** von Strukturknoten  $v$  zu Strukturknoten  $v'$
- **Inhaltlicher Zweifachsprung:** von Strukturknoten  $v$  über Wortknoten zu Strukturknoten  $v'$ 
  - Finde alle Knoten  $v'$ , die über Wortknoten erreichbar sind (Pfadlänge 2)
  - Nehme Top- $q$ -Knoten (Anzahl der Pfade)
  - Wähle zufällig einen davon

- **Struktursprung:** von Strukturknoten  $v$  zu Strukturknoten  $v'$
- **Inhaltlicher Zweifachsprung:** von Strukturknoten  $v$  über Wortknoten zu Strukturknoten  $v'$ 
  - Finde alle Knoten  $v'$ , die über Wortknoten erreichbar sind (Pfadlänge 2)
  - Nehme Top- $q$ -Knoten (Anzahl der Pfade)
  - Wähle zufällig einen davon



- Random Walk
- Gini-Koeffizient
- Inhaltlicher Zweifachsprung

- Random Walk
- Gini-Koeffizient
- Inhaltlicher Zweifachsprung

- Random Walk
- Gini-Koeffizient
- Inhaltlicher Zweifachsprung

- DYCOS ist nur von der lokalen Situation abhängig
- Klassifizierung von einzelnen Knoten möglich
- Klassifizierung ist einfach

- DYCOS ist nur von der lokalen Situation abhängig
- Klassifizierung von einzelnen Knoten möglich
- Klassifizierung ist einfach

- DYCOS ist nur von der lokalen Situation abhängig
- Klassifizierung von einzelnen Knoten möglich
- Klassifizierung ist einfach

Alle folgenden Daten sind der Analyse von Aggarwall und Li entnommen.

Name	Knoten	davon beschriftet	Kanten	Beschriftungen
<b>CORA</b>	19 396	14 814	75 021	5
<b>DBLP</b>	806 635	18 999	4 414 135	5

## ■ Performance:

- Klassifizierung aller Knoten
- Intel Xeon 2.5 GHz mit 32 GB RAM, 1 Kern
- DBLP: < 25 s
- CORA: < 5 s

## ■ Klassifikationsgüte:

- CORA: 82% - 84%
- DBLP: 61% - 66%



## ■ Performance:

### ■ Klassifizierung aller Knoten

- Intel Xeon 2.5 GHz mit 32 GB RAM, 1 Kern
- DBLP: < 25 s
- CORA: < 5 s

## ■ Klassifikationsgüte:

- CORA: 82% - 84%
- DBLP: 61% - 66%

## ■ Performance:

- Klassifizierung aller Knoten
- Intel Xeon 2.5 GHz mit 32 GB RAM, 1 Kern
- DBLP: < 25 s
- CORA: < 5 s

## ■ Klassifikationsgüte:

- CORA: 82% - 84%
- DBLP: 61% - 66%

## ■ Performance:

- Klassifizierung aller Knoten
- Intel Xeon 2.5 GHz mit 32 GB RAM, 1 Kern
- DBLP: < 25 s
- CORA: < 5 s

## ■ Klassifikationsgüte:

- CORA: 82% - 84%
- DBLP: 61% - 66%

## ■ Performance:

- Klassifizierung aller Knoten
- Intel Xeon 2.5 GHz mit 32 GB RAM, 1 Kern
- DBLP: < 25 s
- CORA: < 5 s

## ■ Klassifikationsgüte:

- CORA: 82% - 84%
- DBLP: 61% - 66%

- Performance:
  - Klassifizierung aller Knoten
  - Intel Xeon 2.5 GHz mit 32 GB RAM, 1 Kern
  - DBLP: < 25 s
  - CORA: < 5 s
- Klassifikationsgüte:
  - CORA: 82% - 84%
  - DBLP: 61% - 66%

- Performance:
  - Klassifizierung aller Knoten
  - Intel Xeon 2.5 GHz mit 32 GB RAM, 1 Kern
  - DBLP: < 25 s
  - CORA: < 5 s
- Klassifikationsgüte:
  - CORA: 82% - 84%
  - DBLP: 61% - 66%

- Performance:
  - Klassifizierung aller Knoten
  - Intel Xeon 2.5 GHz mit 32 GB RAM, 1 Kern
  - DBLP: < 25 s
  - CORA: < 5 s
- Klassifikationsgüte:
  - CORA: 82% - 84%
  - DBLP: 61% - 66%

# Danke!

## Gibt es Fragen?



- Crystal\_Clear\_app\_personal.png von [Wikipedia Commons](#)

- Charu C. Aggarwal, Nan Li: *On Node Classification in Dynamic Content-based Networks*.
- Smriti Bhagat, Graham Cormode und S. Muthukrishnan. *Node Classification in Social Networks*.
- M. F. Porter. Readings in Information Retrieval. Kapitel *An Algorithm for Suffix Stripping*.
- Jeffrey S. Vitter. *Random Sampling with a Reservoir*.

Der Foliensatz und die  $\text{\LaTeX}$  und TikZ-Quellen sind unter  
[github.com/MartinThoma/LaTeX-examples/tree/master/presentations/Datamining-Proseminar](https://github.com/MartinThoma/LaTeX-examples/tree/master/presentations/Datamining-Proseminar)  
Kurz-URL: [tinyurl.com/Info-Proseminar](https://tinyurl.com/Info-Proseminar)