

Subsampling of Complex Networks



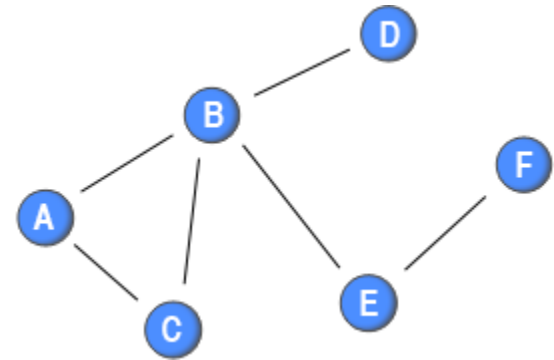
TECHNISCHE
UNIVERSITÄT
DARMSTADT

Kai Rathmann

Supervisor: Prof. Dr. Thorsten Strufe

Coordinator: Benjamin Schiller

- Complex Network $G = (V, E)$
 - Menge von Knoten V
 - Menge von Kanten $E \subset (V \times V)$
 - keine Schleifen
 - zusammenhängend, endlich
- Subsampling
 - Auswahl einer Teilmenge von Knoten $V' \subset V$ und Kanten $E' \subset E$



- Problemstellung: Großes Netzwerk ist nicht vollständig verwendbar
 - Ressourcenaufwand
 - Zugriffsbeschränkungen
- Gewinn eines repräsentativen Teilnetzwerkes durch Subsampling
- Frage: Welche Eigenschaften des Gesamtnetzwerkes werden durch das Subsampling wie erhalten?

Vorgehen

- Gruppierung vieler bekannter Subsampling-Algorithmen
- Versuch einer Klassifikation
- Skizzierung neuer Varianten
- Implementierung in GTNA
- Evaluation

- Metriken
- Subsampling-Algorithmen
- Implementierung
- Evaluation
- Zusammenfassung und Fazit

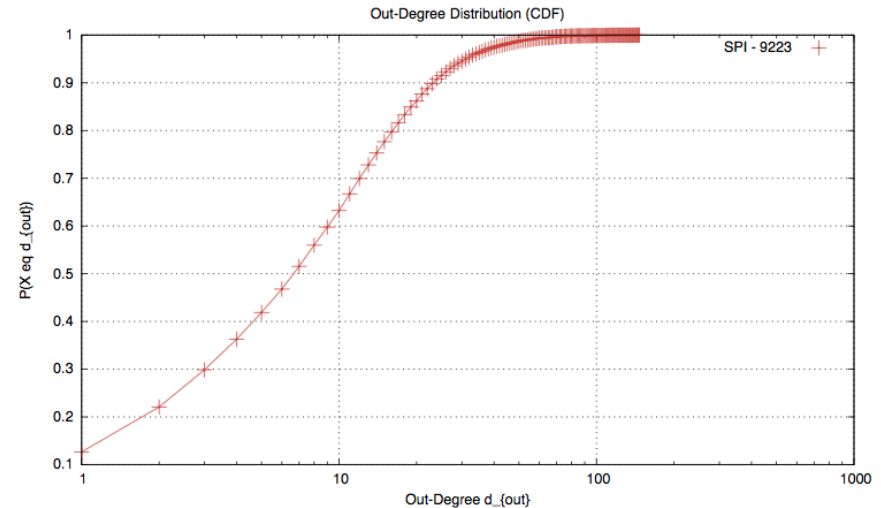


- Gradverteilung
- Durchmesser
- Charakteristische Pfadlänge
- Clusterkoeffizient
- Netzwerkgröße

Metriken: Gradverteilung



- Grad k_v eines Knoten v : Anzahl Kanten, die inzident zu diesem Knoten sind
 - Ausgangsgrad k_v^+ : ausgehende Kanten (v, w)
 - Eingangsgrad k_v^- : eingehende Kanten (w, v)
- Gradverteilung: Verteilung von Knotengraden
- Skalenfreie Netzwerke
 - $P(k) \sim k^{-\gamma}$



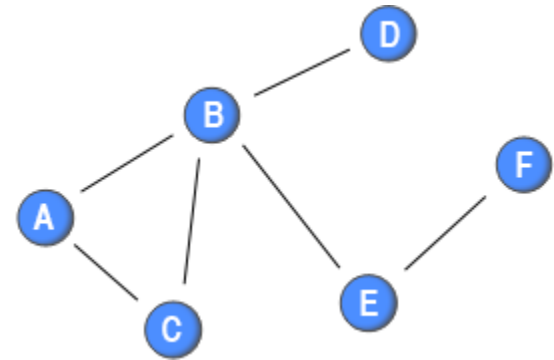
Metriken: Durchmesser, Char. Pfadlänge



- Pfad (v_1, \dots, v_n) : Knoten v_1, \dots, v_n , so dass Kanten $(v_1, v_2), \dots, (v_{n-1}, v_n)$ existieren
- Länge eines Pfades (v_1, \dots, v_n) : $n - 1$
- Abstand zweier Knoten: Länge des kürzesten Pfades
- Durchmesser: Maximaler Abstand zweier Knoten im Netzwerk
- Charakteristische Pfadlänge: Durchschnittlicher Abstand zweier Knoten im Netzwerk
- Milgram's Experiment
- Six Degrees of Separation

Metriken: Clusterkoeffizient

- Triplet:
 - drei Knoten, die durch zwei (offenes Triplet) oder drei (geschlossenes Triplet) Kanten verbunden sind
- Clusterkoeffizient:
 - Anzahl geschlossener Triplets / Anzahl aller Triplets



Metriken: Netzwerkgröße

- Anzahl Knoten im Netzwerk
- relevant nur für *wiederkehrende* Subsampling-Algorithmen

Subsampling-Algorithmen

- Uniform Node Sampling
- Classic Random Walk
- Metropolis-Hastings Random Walk
- Frontier Sampling
- Breadth-First Search
- Depth-First Search
- Forest Fire
- Snowball Sampling
- Respondent-Driven Sampling

Algorithmen: Uniform Node Sampling



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Uniforme Auswahl von Knoten

Algorithmen: Classic Random Walk

- Uniforme Wahl des Startknotens
- In jedem Schritt
 - Uniforme Wahl eines Nachbarknotens
 - $p(v, w) = 1/k_v$
- Wiederholte Wahl von Knoten möglich

Algorithmen: Metropolis-Hastings R.W.



- Uniforme Wahl des Startknotens
- In jedem Schritt
 - Gewichtete Wahl eines Nachbarknotens
 - $p(v, w) = 1/k_v * \min(1, k_v/k_w)$ falls w Nachbar von v
 - $p(v, v) = 1 - \sum_{w \neq v} p(v, w)$
- Wiederholte Wahl von Knoten möglich

Algorithmen: Frontier Sampling

- m Instanzen eines klassischen Random Walks
- In jedem Schritt wähle einen Random Walk aus, um einen Schritt auszuführen
 - $p(i) = k(v_i) / \sum_i k(v_i)$

Algorithmen: Breadth-First Search

- Markiere den Startknoten
- Speichere den Startknoten in der Schlange
- Solange die Schlange nicht leer ist
 - Entnimm den ersten Knoten v aus der Schlange
 - Für alle (ausgehenden) Kanten (v, w)
 - Falls w nicht markiert ist
 - Markiere w
 - Speichere w in der Schlange

Algorithmen: Depth-First Search



- Markiere v
- Für alle (ausgehenden) Kanten (v, w)
 - Falls w unmarkiert ist
 - Depth-First Search(w)

Algorithmen: Forest Fire

- Ähnlich Breadth-First Search
- „Markiere w “ für alle unmarkierten Nachbarknoten nur mit Wahrscheinlichkeit p
 - $p=1$: BFS

Algorithmen: Snowball Sampling

- Markiere den Startknoten
- Speichere den Startknoten in der Schlange
- Solange die Schlange nicht leer ist
 - Entnimm den ersten Knoten v aus der Schlange
 - Für max. n (ausgehende) Kanten (v, w) mit unmarkiertem w
 - Markiere w
 - Speichere w in der Schlange

Algorithmen: Respondent-Driven Sampling



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Ähnlich Snowball Sampling
- Wiederholte Wahl von Knoten möglich

Algorithmen: Klassifikation



| Algorithmus | Nachbarn |
|-------------|--------------|
| UNS | keine |
| RW | 1 |
| MHRW | 1 |
| FS | 1 |
| BFS | alle |
| DFS | 1 |
| FF | alle (mit p) |
| SS | mehrere |
| RDS | mehrere |

Algorithmen: Klassifikation



| Nachbarn | Algorithmus |
|------------------|-------------------|
| $n = 0$ | UNS |
| $n = 1$ | RW, MHRW, FS, DFS |
| $1 < n < \infty$ | SS, RDS |
| $n = \infty$ | BFS, FF |

Algorithmen: Klassifikation



| Nachbarn | nicht wiederkehrend | wiederkehrend |
|------------------|---------------------|---------------|
| $n = 0$ | UNS | |
| $n = 1$ | DFS | RW, MHRW, FS |
| $1 < n < \infty$ | SS | RDS |
| $n = \infty$ | BFS, FF | |

Algorithmen: Klassifikation



| Nachbarn | | nicht wiederkehrend | wiederkehrend |
|------------------|-----------|---------------------|---------------|
| $n = 0$ | | UNS | |
| $n = 1$ | geordnet | DFS | |
| | uniform | | RW, FS |
| | gewichtet | | MHRW |
| $1 < n < \infty$ | uniform | SS | RDS |
| $n = \infty$ | alle | BFS | |
| | p fest | FF | |

Algorithmen: Klassifikation



| Nachbarn | | nicht wiederkehrend | wiederkehrend | mehr-dimensional |
|------------------|-----------|---------------------|---------------|------------------|
| $n = 0$ | | UNS | | |
| $n = 1$ | geordnet | DFS | | |
| | uniform | | RW | FS |
| | gewichtet | | MHRW | |
| $1 < n < \infty$ | uniform | SS | RDS | |
| $n = \infty$ | alle | BFS | | |
| | p fest | FF | | |

Algorithmen: Klassifikation



| Nachbarn | | nicht wiederkehrend | wiederkehrend | unabhängig mehrdimensional | abhängig mehrdimensional |
|------------------|-------------|---------------------|---------------|----------------------------|--------------------------|
| $n = 0$ | | UNS | | | |
| $n = 1$ | geordnet | DFS | | | |
| | uniform | | RW | | FS |
| | gewichtet | | MHRW | | |
| $1 < n < \infty$ | geordnet | | | | |
| | uniform | SS | RDS | | |
| | gewichtet | | | | |
| $n = \infty$ | alle | BFS | | | |
| | p fest | FF | | | |
| | p gewichtet | | | | |

Algorithmen: Klassifikation



| Nachbarn | | nicht wiederkehrend | wiederkehrend | unabhängig mehrdimensional | abhängig mehrdimensional |
|------------------|-------------|---------------------|---------------|----------------------------|--------------------------|
| $n = 0$ | | UNS | | | |
| $n = 1$ | geordnet | DFS | | | |
| | uniform | * | RW | * | FS |
| | gewichtet | * | MHRW | * | * |
| $1 < n < \infty$ | geordnet | * | | | |
| | uniform | SS | RDS | * | * |
| | gewichtet | * | * | * | * |
| $n = \infty$ | alle | BFS | | | |
| | p fest | FF | | | |
| | p gewichtet | * | | | |

Algorithmen: Klassifikation



| Nachbarn | | nicht wiederkehrend | wiederkehrend | unabhängig mehr-dimensional | abhängig mehr-dimensional |
|------------------|-------------|--------------------------|--------------------------|-----------------------------|---------------------------|
| $n = 0$ | | UNS | | | |
| $n = 1$ | geordnet | DFS_{geordnet} | | | |
| | uniform | DFS_{uniform} | RW_{uniform} | mRW_{uniform} | FS_{uniform} |
| | gewichtet | $DFS_{\text{gewichtet}}$ | $RW_{\text{gewichtet}}$ | $mRW_{\text{gewichtet}}$ | $FS_{\text{gewichtet}}$ |
| $1 < n < \infty$ | geordnet | SS_{geordnet} | | | |
| | uniform | SS_{uniform} | RDS_{uniform} | $mRDS_{\text{uniform}}$ | $FS+_{\text{uniform}}$ |
| | gewichtet | $SS_{\text{gewichtet}}$ | $RDS_{\text{gewichtet}}$ | $mRDS_{\text{gewicht}}$ | $FS+_{\text{gewichtet}}$ |
| $n = \infty$ | alle | BFS | | | |
| | p fest | FF_{fest} | | | |
| | p gewichtet | $FF_{\text{gewichtet}}$ | | | |

- Transformation
 - `isApplicable(Graph)`:
boolean
 - `transform(Graph)`: Graph
- Subsampling
 - `subsample(Graph)`:
Collection<Node>
 - Normalisierung
- Subklassen für alle Algorithmen



GTNA
Graph-Theoretic Network Analyzer

Evaluation

- Setup
- Netzwerke
- Ergebnisse

Evaluation: Setup

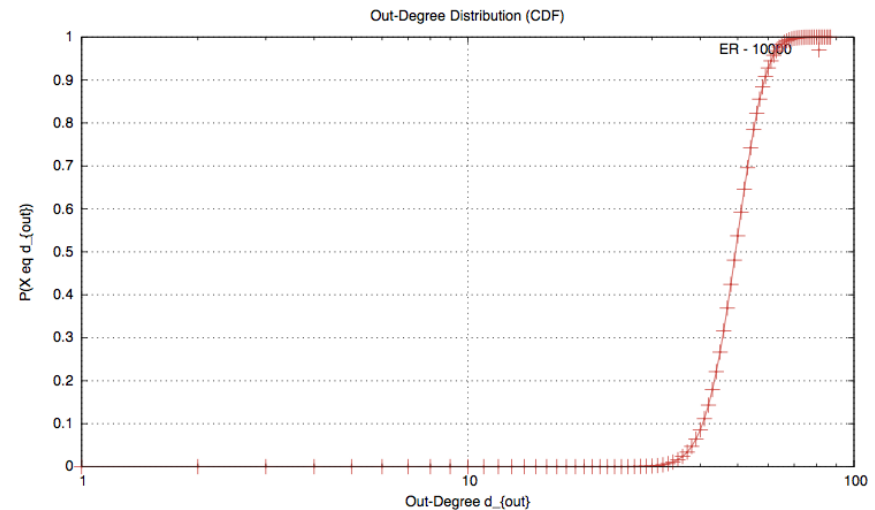
- 100 Durchläufe für jede Kombination aus Netzwerktyp und Subsampling-Algorithmus
 - Netzwerkmodelle: 1 Durchlauf pro generiertem Netzwerk
 - Reale Netzwerke: 100 Durchläufe
- 10% Budget
- $n = 3$
- $p = 0,7$
- $m = 10$

Evaluation: Netzwerke

- Erdős–Rényi
- Barabási–Albert
- Studentenportal Ilmenau
- Web of Trust (ungerichtet)
- Web of Trust (gerichtet)

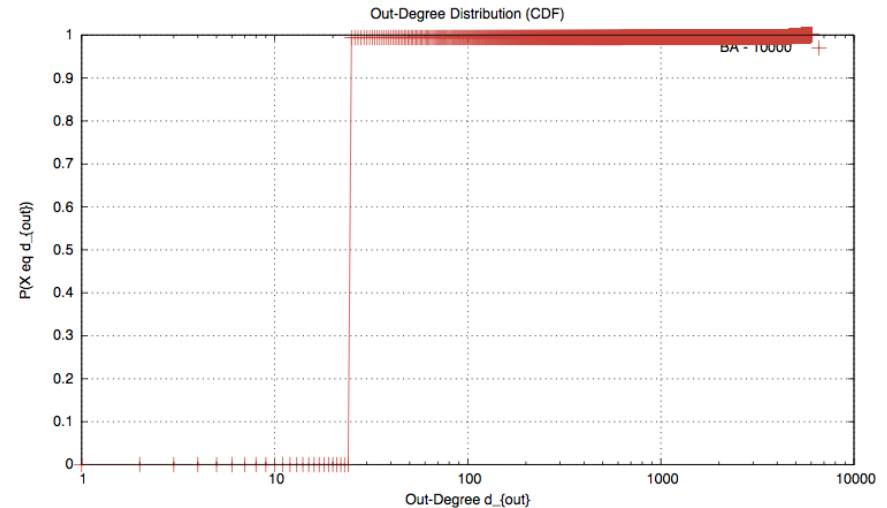
Eval. – Netzwerke: Erdős–Rényi

- Setzt Kanten zwischen Knotenpaaren mit gleicher Wahrscheinlichkeit, unabhängig von anderen Kanten
- 10.000 Knoten
- Durchschnittsgrad 50
- Durchschnittswerte:
 - Durchmesser 4,0
 - Ch. Pfadlänge 2,77
 - Clusterkoeffizient 0.0050



Eval. – Netzwerke: Barabási–Albert

- Preferential attachment
- Skalenfreies Netzwerk
- 10.000 Knoten
- Kantenzahl pro Knoten 25
- Durchschnittswerte:
 - Durchmesser 2,3
 - Ch. Pfadlänge 2,00
 - Clusterkoeffizient 0,4592

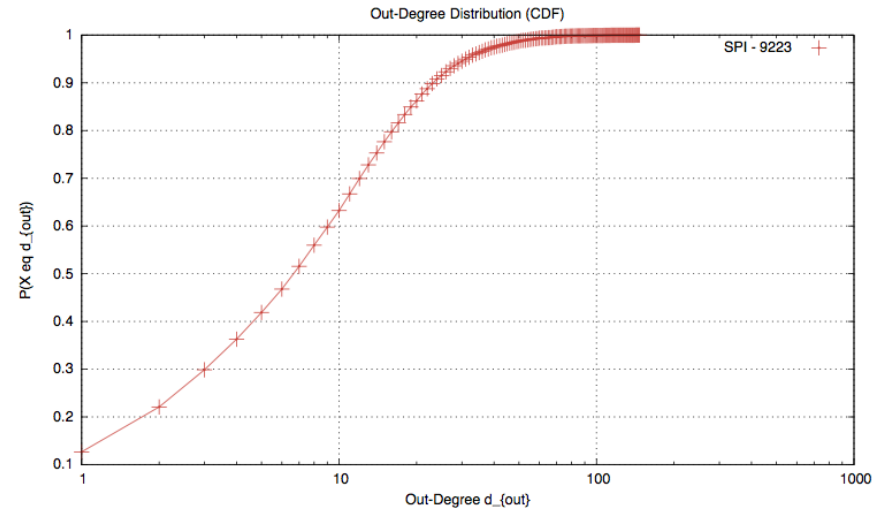


Eval. – Netzwerke: Studentenportal Ilmenau



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Soziales Netzwerk
- 9.223 Knoten
- Durchmesser 12
- Ch. Pfadlänge 4,66
- Clusterkoeffizient 0,2964

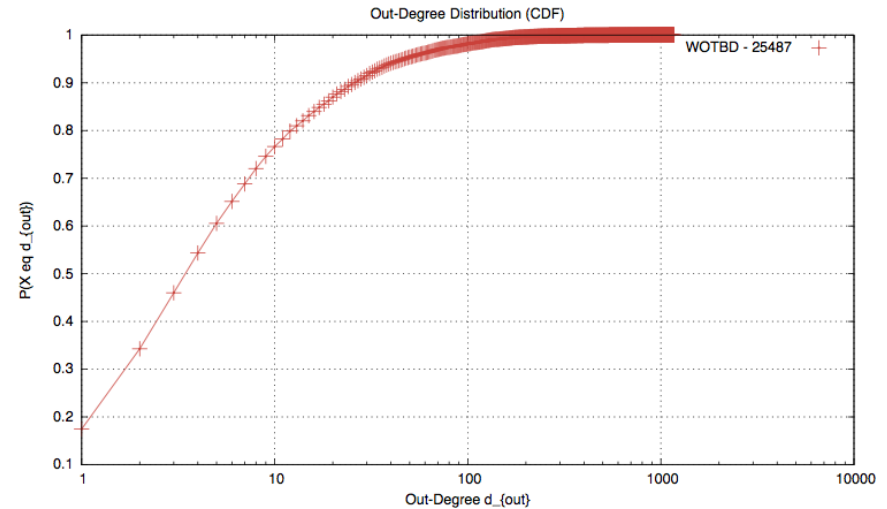


Eval. – Netzwerke: Web of Trust

- Zertifizierungsgraph von PGP
- Evaluation sowohl der ungerichteten als auch der gerichteten Variante
- 25.487 Knoten

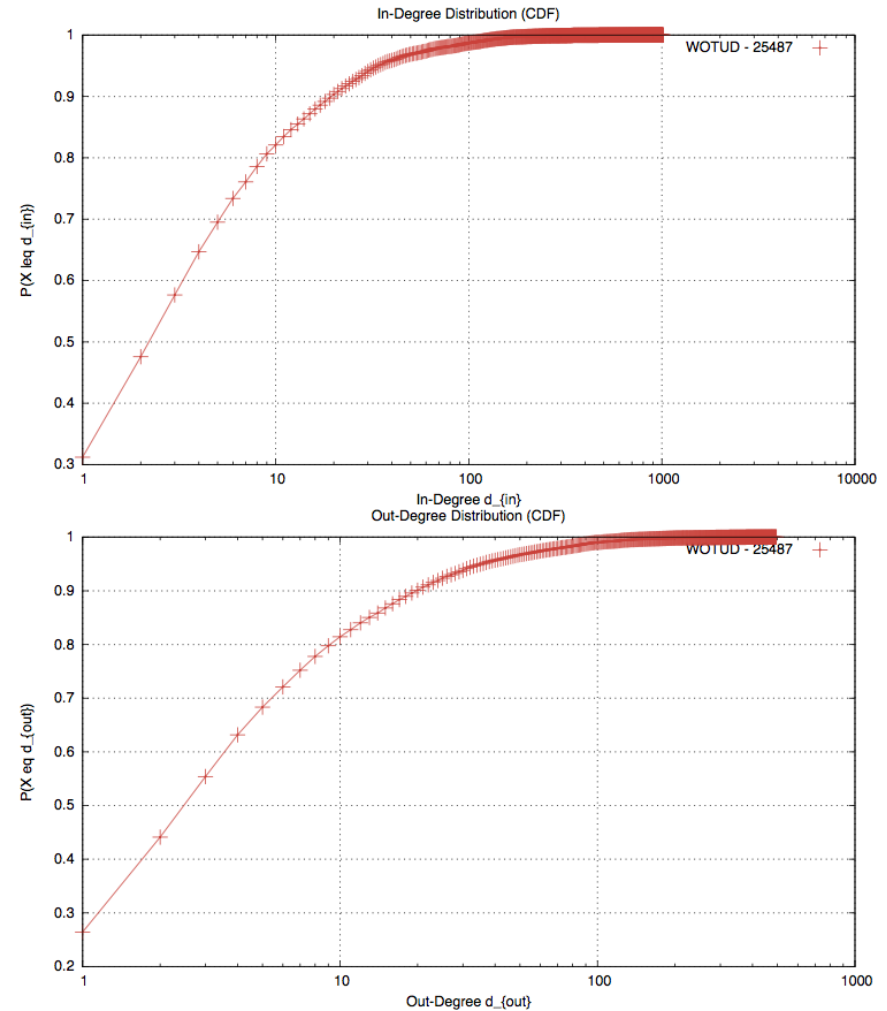
Eval. – Netzwerke: Web of Trust (ungerichtet)

- Durchmesser 15
- Ch. Pfadlänge 5,07
- Clusterkoeffizient 0,4665



Eval. – Netzwerke: Web of Trust (gerichtet)

- Durchmesser 25
- Ch. Pfadlänge 5,99
- Clusterkoeffizient 0,3707



Evaluation: Ergebnisse

- Netzwerkgröße
- Gradverteilung
- Durchmesser
- Charakteristische Pfadlänge
- Clusterkoeffizient

Eval. – Ergebnisse: Netzwerkgröße



| Nachbarn | | nicht wiederkehrend | wiederkehrend | unabhängig mehrdimensional | abhängig mehrdimensional |
|------------------|-------------|--------------------------|--------------------------|----------------------------|--------------------------|
| $n = 0$ | | UNS | | | |
| $n = 1$ | geordnet | DFS_{geordnet} | | | |
| | uniform | DFS_{uniform} | RW_{uniform} | mRW_{uniform} | FS_{uniform} |
| | gewichtet | $DFS_{\text{gewichtet}}$ | $RW_{\text{gewichtet}}$ | $mRW_{\text{gewichtet}}$ | $FS_{\text{gewichtet}}$ |
| $1 < n < \infty$ | geordnet | SS_{geordnet} | | | |
| | uniform | SS_{uniform} | RDS_{uniform} | $mRDS_{\text{uniform}}$ | $FS+_{\text{uniform}}$ |
| | gewichtet | $SS_{\text{gewichtet}}$ | $RDS_{\text{gewichtet}}$ | $mRDS_{\text{gewichtet}}$ | $FS+_{\text{gewichtet}}$ |
| $n = \infty$ | alle | BFS | | | |
| | p fest | FF_{fest} | | | |
| | p gewichtet | $FF_{\text{gewichtet}}$ | | | |

Eval. – Ergebnisse: Netzwerkgröße



| Nachbarn | | nicht wiederkehrend | wiederkehrend | unabhängig mehrdimensional | abhängig mehrdimensional |
|------------------|-------------|---------------------|--------------------------|----------------------------|--------------------------|
| $n = 0$ | | | | | |
| $n = 1$ | geordnet | | | | |
| | uniform | | RW_{uniform} | mRW_{uniform} | FS_{uniform} |
| | gewichtet | | $RW_{\text{gewichtet}}$ | $mRW_{\text{gewichtet}}$ | $FS_{\text{gewichtet}}$ |
| $1 < n < \infty$ | geordnet | | | | |
| | uniform | | RDS_{uniform} | $mRDS_{\text{uniform}}$ | $FS+_{\text{uniform}}$ |
| | gewichtet | | $RDS_{\text{gewichtet}}$ | $mRDS_{\text{gewichtet}}$ | $FS+_{\text{gewichtet}}$ |
| $n = \infty$ | alle | | | | |
| | p fest | | | | |
| | p gewichtet | | | | |

- Reale Netzwerke
 - Primärer Einfluss: Nachbarstrategie
 1. $n = 1$ uniform
 2. $1 < n < \infty$ gewichtet
 3. $1 < n < \infty$ uniform
 4. $n = 1$ gewichtet
 - Sekundärer Einfluss: Dimensionale Strategie
 1. abhängig mehrdimensional
 2. unabhängig mehrdimensional
 3. eindimensional

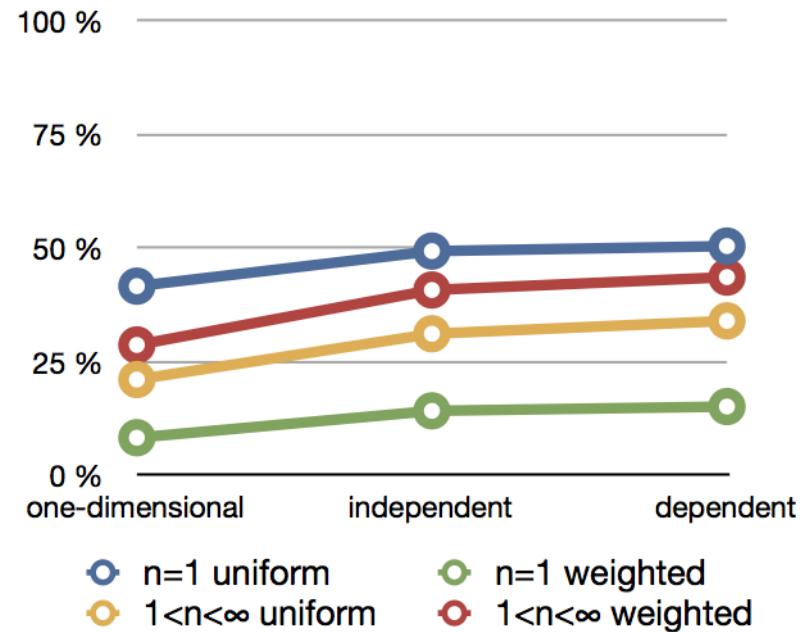
Eval. – Ergebnisse: Netzwerkgröße



| Nachbarn | | nicht wiederkehrend | wiederkehrend | unabhängig mehr-dimensional | abhängig mehr-dimensional |
|------------------|-------------|---------------------|--------------------------|-----------------------------|---------------------------|
| $n = 0$ | | | | | |
| $n = 1$ | geordnet | | | | |
| | uniform | | RW_{uniform} | mRW_{uniform} | FS_{uniform} |
| | gewichtet | | $RW_{\text{gewichtet}}$ | $mRW_{\text{gewichtet}}$ | $FS_{\text{gewichtet}}$ |
| $1 < n < \infty$ | geordnet | | | | |
| | uniform | | RDS_{uniform} | $mRDS_{\text{uniform}}$ | $FS+_{\text{uniform}}$ |
| | gewichtet | | $RDS_{\text{gewichtet}}$ | $mRDS_{\text{gewichtet}}$ | $FS+_{\text{gewichtet}}$ |
| $n = \infty$ | alle | | | | |
| | p fest | | | | |
| | p gewichtet | | | | |

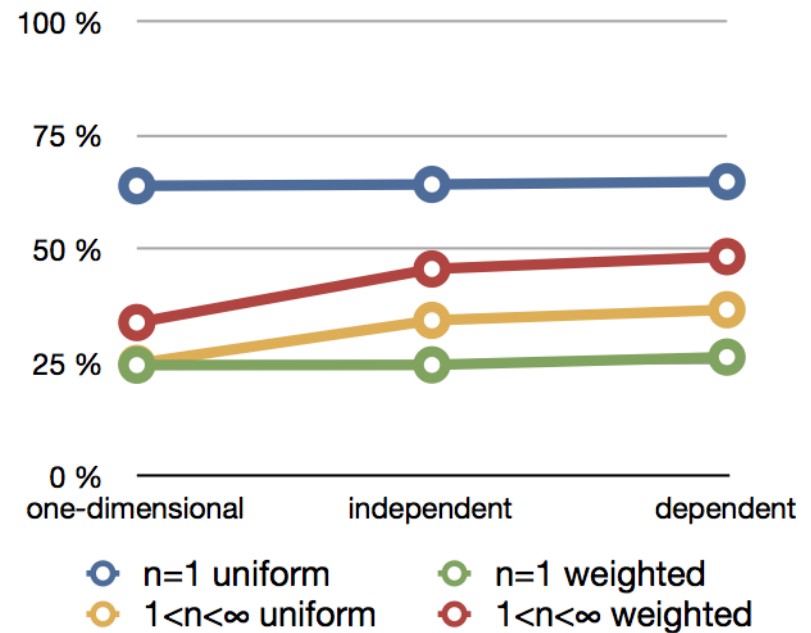
Eval. – Ergebnisse: Netzwerkgröße

Web of Trust (gerichtet)



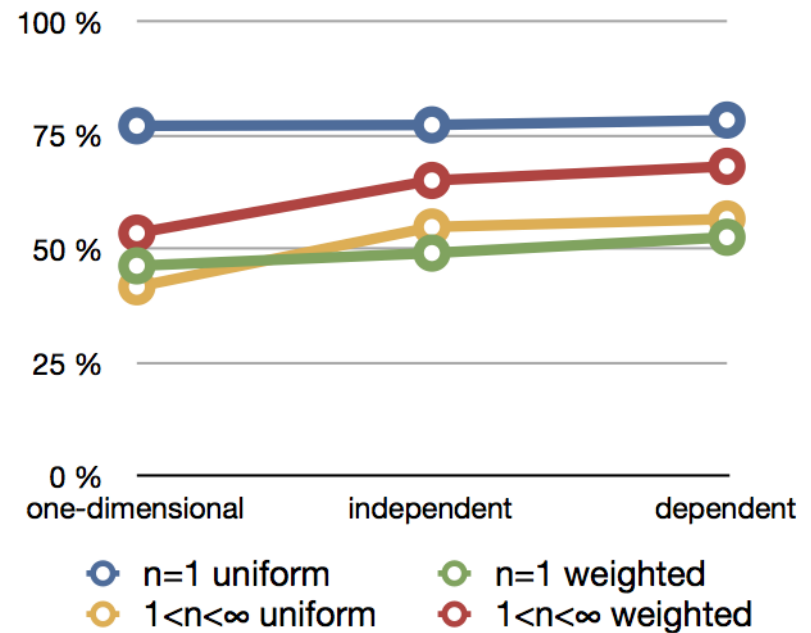
Eval. – Ergebnisse: Netzwerkgröße

Web of Trust (ungerichtet)



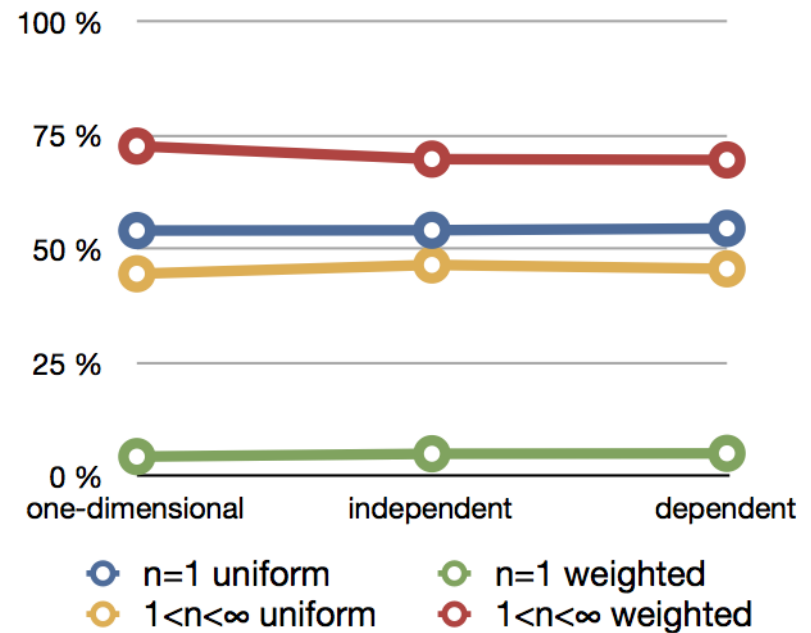
Eval. – Ergebnisse: Netzwerkgröße

Studentenportal Ilmenau



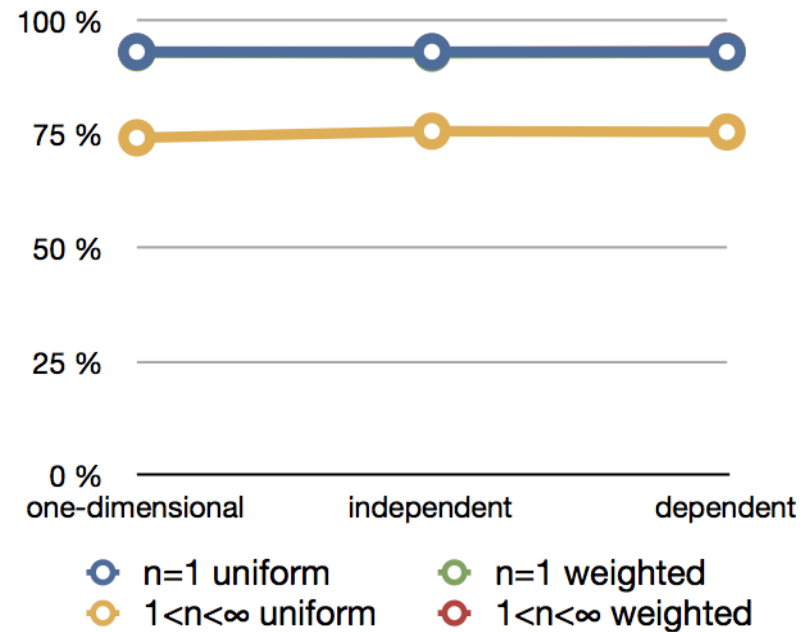
Eval. – Ergebnisse: Netzwerkgröße

Barabási–Albert



Eval. – Ergebnisse: Netzwerkgröße

Erdős–Rényi





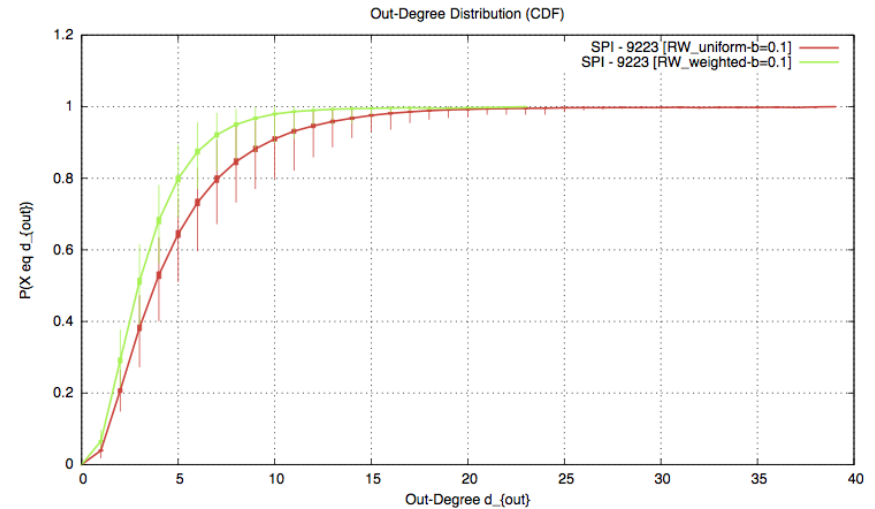
Eval. – Ergebnisse: Gradverteilung

- Ergebnisse zu $RW_{\text{gewichtet}}$ und FS_{uniform} reproduziert
- Effekt von Gewichtung gezeigt für SS, mRW, FS, RDS, mRDS, FS+, FF
- $DFS_{\text{uniform}} \approx DFS_{\text{gewichtet}}$
- Konzeptkombination $FS_{\text{gewichtet}}$ verstärkt Effekt

Eval. – Ergebnisse: Gradverteilung

RW_{uniform} vs. $RW_{\text{gewichtet}}$

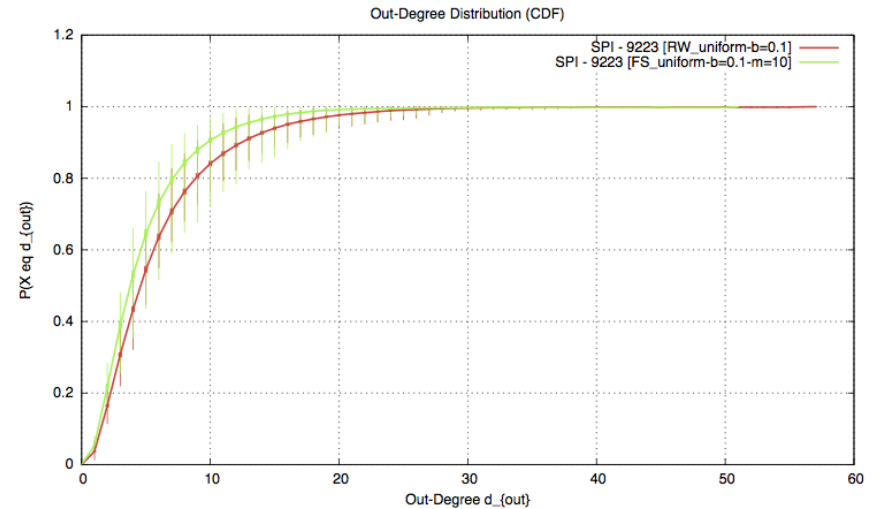
am Beispiel Studentenportal Ilmenau



Eval. – Ergebnisse: Gradverteilung

RW_{uniform} vs. FS_{uniform}

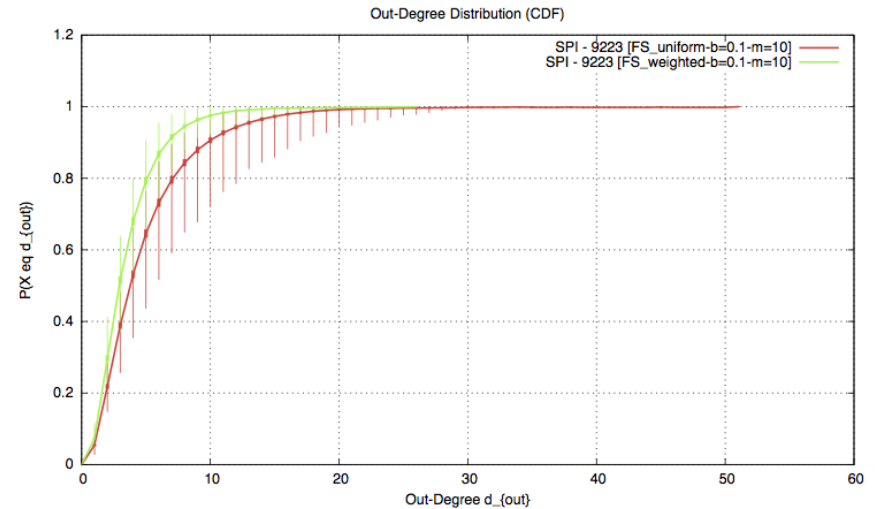
am Beispiel Studentenportal Ilmenau



Eval. – Ergebnisse: Gradverteilung

FS_{uniform} vs. $FS_{\text{gewichtet}}$

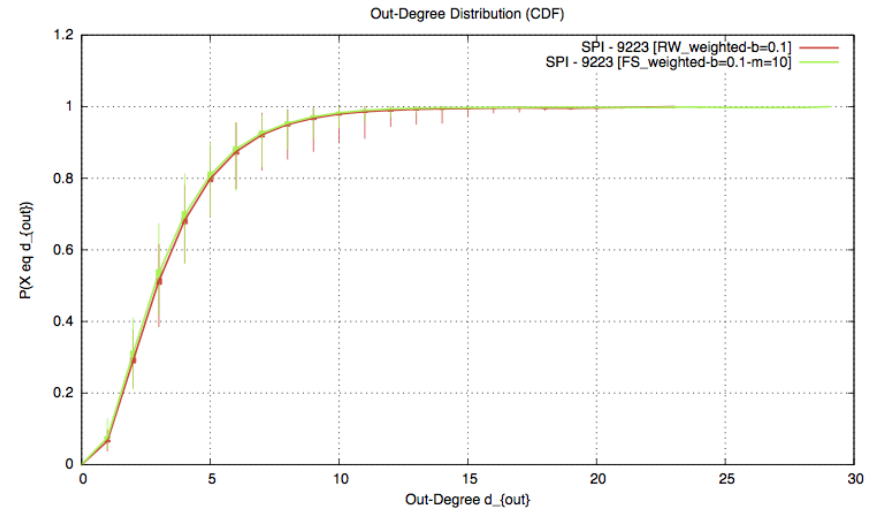
am Beispiel Studentenportal Ilmenau



Eval. – Ergebnisse: Gradverteilung

$RW_{\text{gewichtet}}$ vs. $FS_{\text{gewichtet}}$

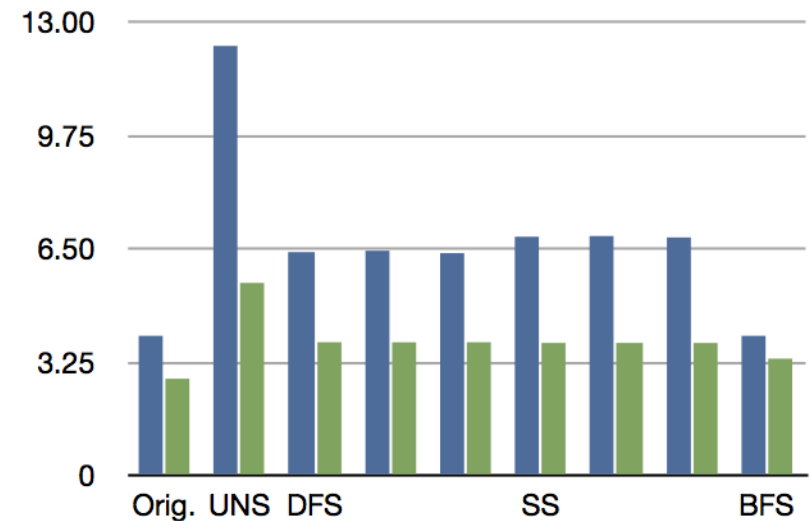
am Beispiel Studentenportal Ilmenau



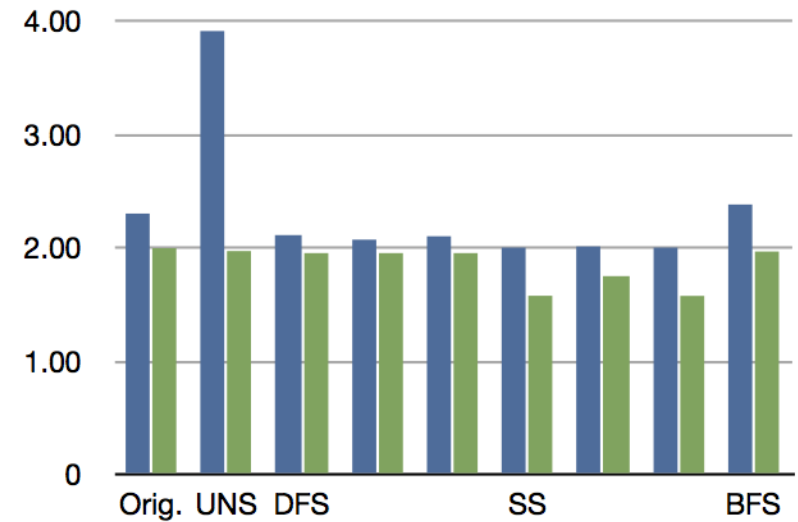


- DFS-Varianten insgesamt am originaltreuesten
- SS-Varianten produzieren niedrigere Werte als DFS-Varianten in allen evaluierten Netzwerken außer Erdős–Rényi
- UNS produziert höchste Werte in den ungerichteten Netzwerken, sehr niedrige in dem gerichteten Netzwerk
- BFS produziert niedrigste Werte in allen evaluierten Netzwerken außer Barabási–Albert

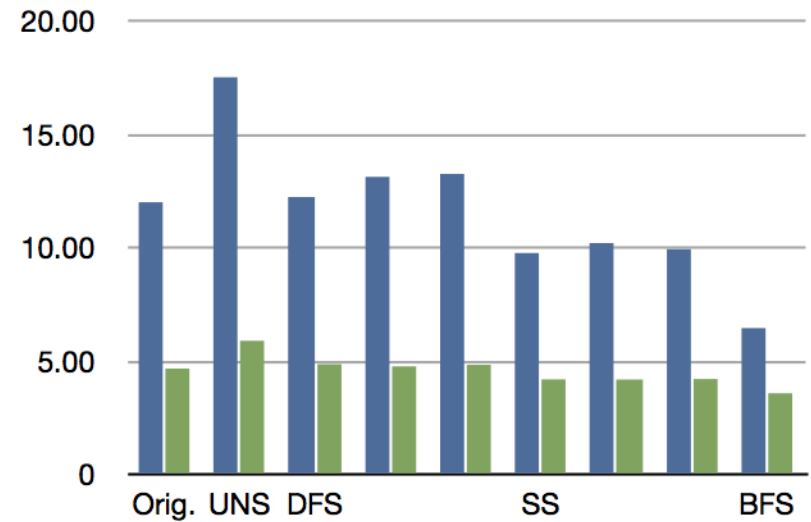
Erdős–Rényi



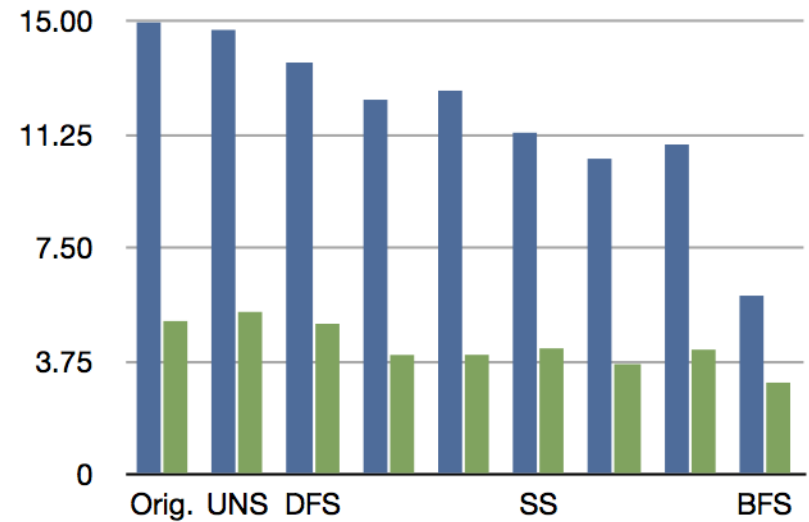
Barabási–Albert



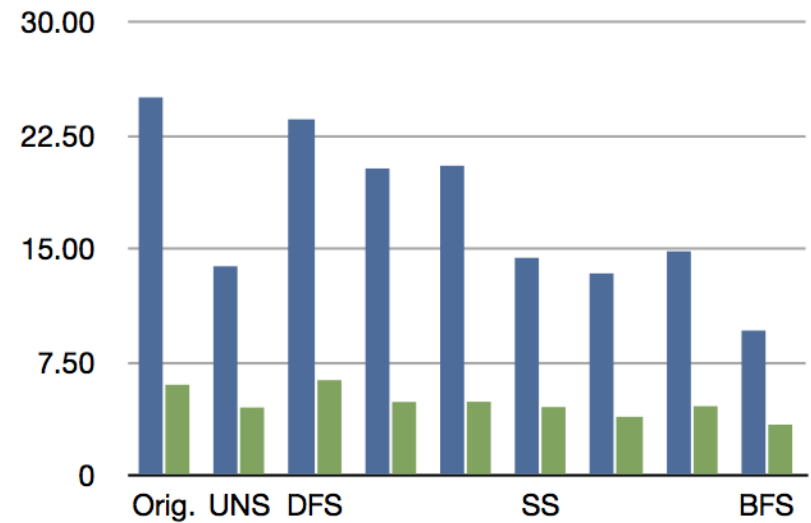
Studentenportal Ilmenau



Web of Trust (ungerichtet)



Web of Trust (gerichtet)



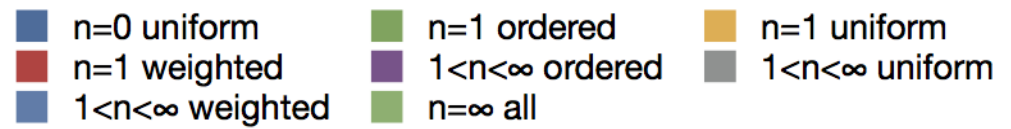
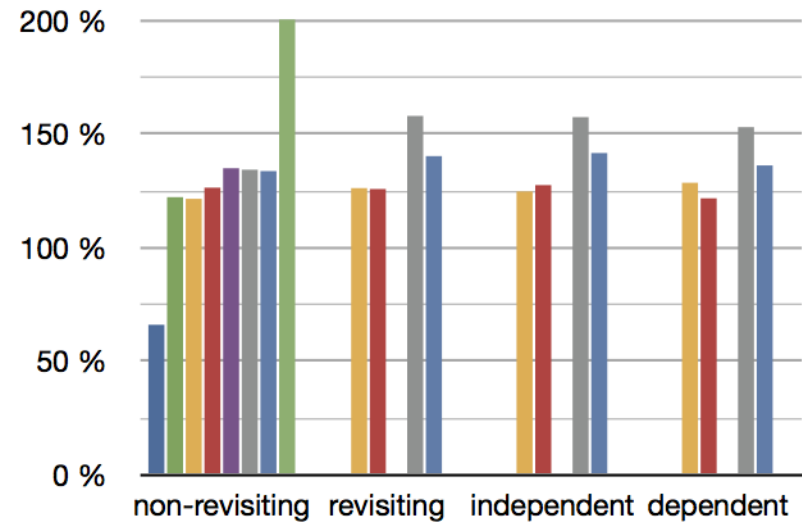
Eval. – Ergebnisse: Clusterkoeffizient



- Algorithmen mit $0 < n < \infty$ produzieren Werte relativ dicht am Originalwert mit geringen Abweichungen
- UNS produziert kleinste Werte (19,6% – 65,9% des Originalwertes)

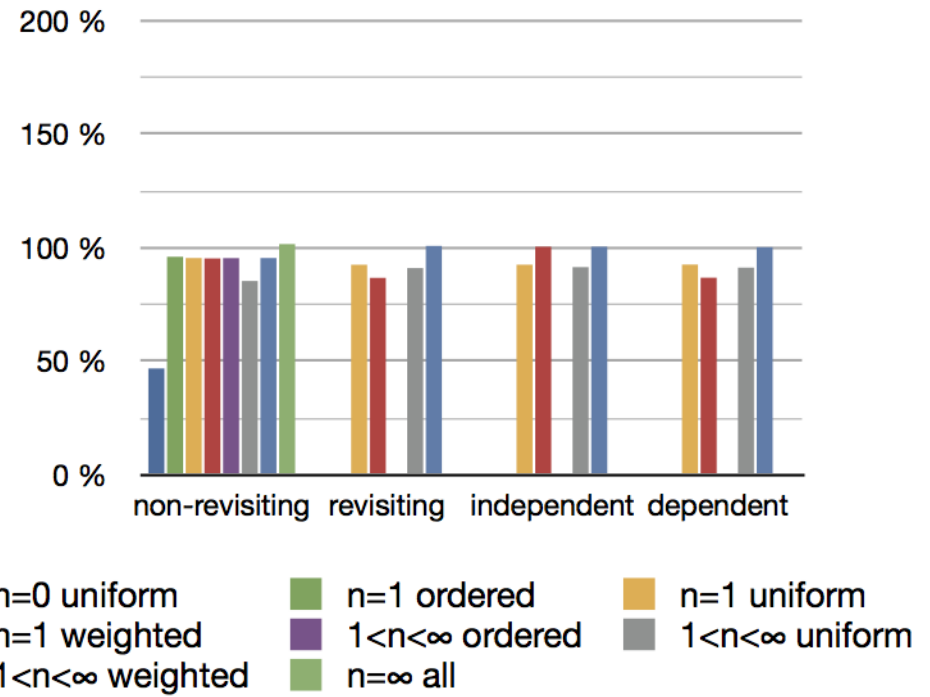
Eval. – Ergebnisse: Clusterkoeffizient

Erdős–Rényi



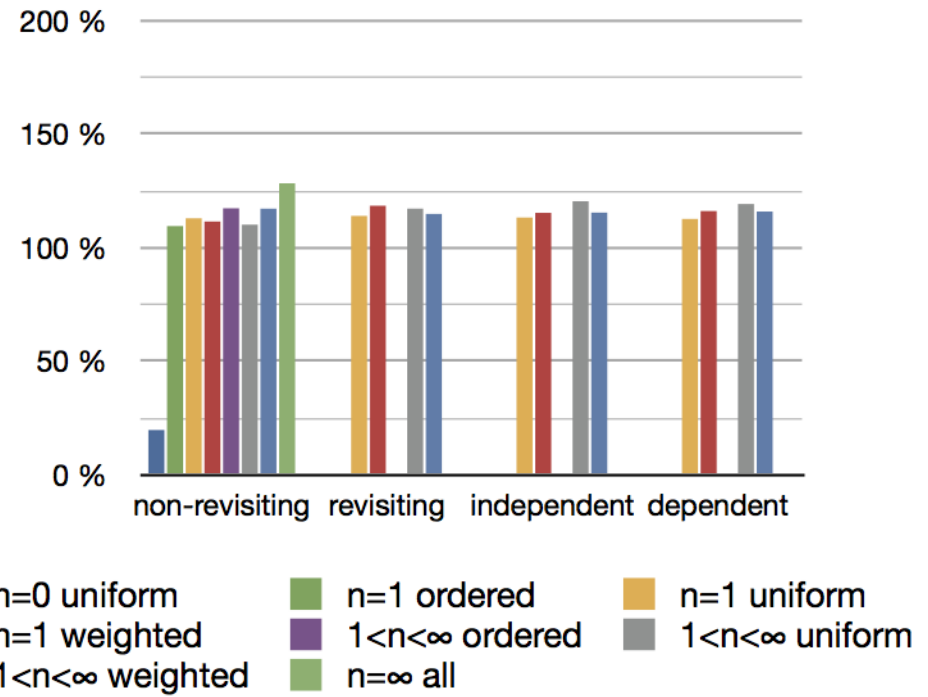
Eval. – Ergebnisse: Clusterkoeffizient

Barabási–Albert



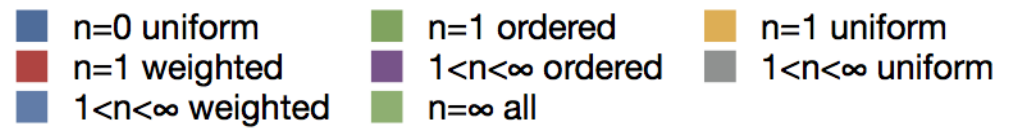
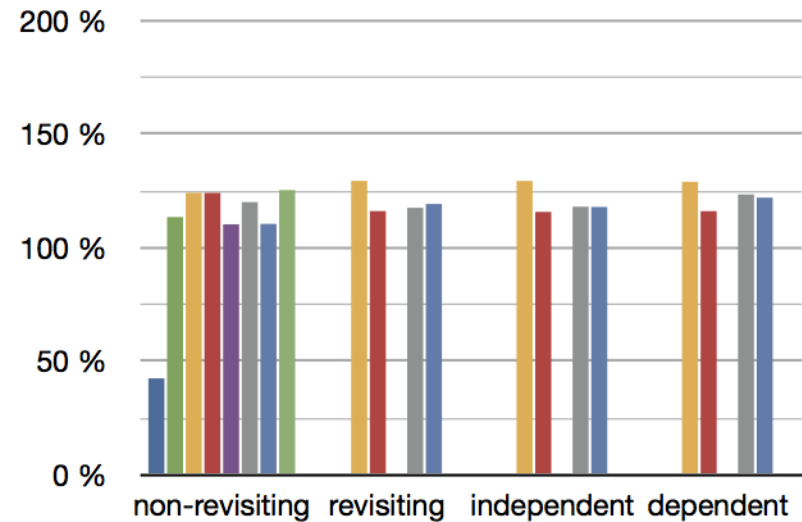
Eval. – Ergebnisse: Clusterkoeffizient

Studentenportal Ilmenau



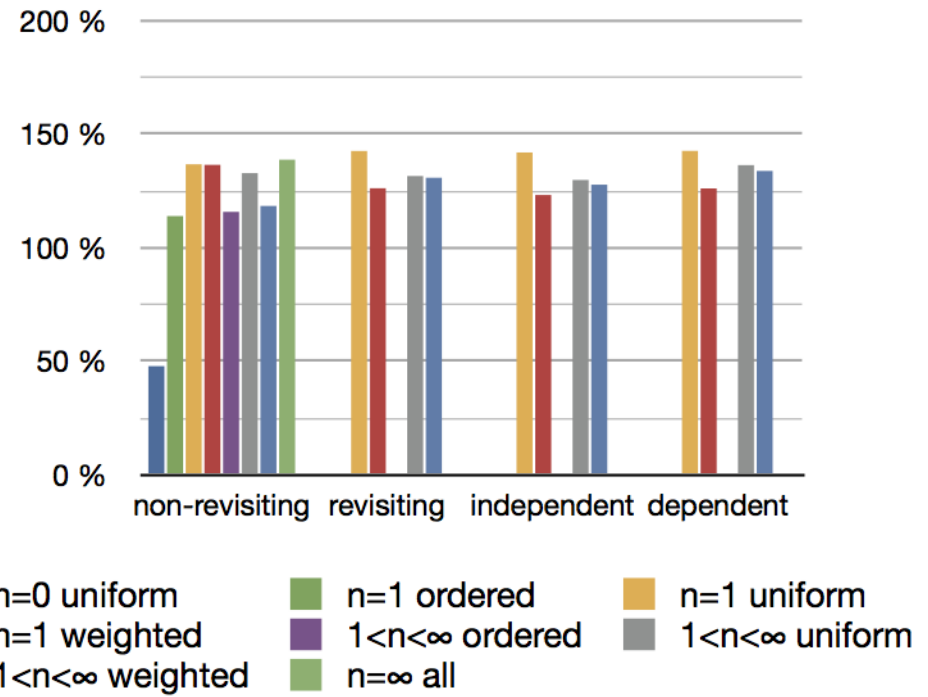
Eval. – Ergebnisse: Clusterkoeffizient

Web of Trust (ungerichtet)



Eval. – Ergebnisse: Clusterkoeffizient

Web of Trust (gerichtet)



Zusammenfassung

- Subsampling-Algorithmen zweidimensional klassifiziert nach Nachbarwahl und dimensionaler Strategie
- Implementierung
- Evaluation

- Netzwerkgröße primär durch Nachbarwahl und sekundär durch dimensionale Strategie beeinflusst
- Kombination von Konzepten bewirkt Kombination bzw. Verstärkung von Ergebnissen der Gradverteilung
- DFS-Varianten am originaltreuesten für Durchmesser und Charakteristische Pfadlänge
- Clusterkoeffizient gut bewahrt für $0 < n < \infty$
- $n = 0$ und $n = \infty$ produzieren extreme Ergebnisse



Subsampling of Complex Networks