Introduction
0000

Green measures
00000

Methods Compared
0000000000000

Experiment on Wikipedia
0000

Conclusion
00
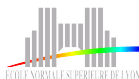
# Finding Related Pages Using Green Measures: The Example of Wikipedia

Yann Ollivier          Pierre Senellart

CENTRE NATIONAL
DE LA RECHERCHE
SCIENTIFIQUE

ÉCOLE NORMALE SUPÉRIEURE DE LYON

INSTITUT NATIONAL
DE RECHERCHE
EN INFORMATIQUE
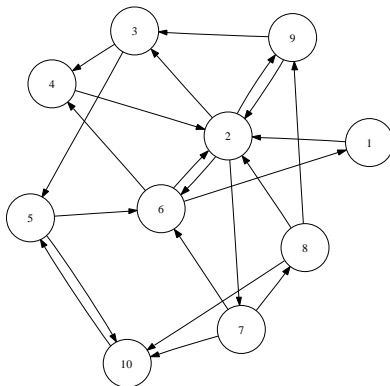ET EN AUTOMATIQUE

$INRIA$
FUTURS

UNIVERSITÉ
PARIS-SUD 11

*AAAI*

July 24th, 2007

# Related nodes in a graph
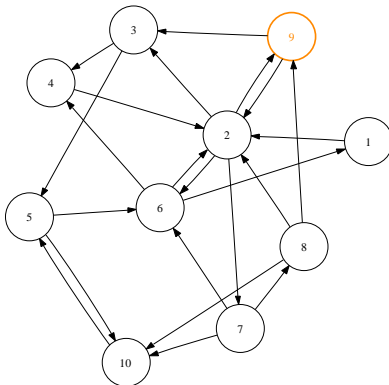
Given a hyperlinked environment (= a graph)...



## Problem

Finding nodes semantically related to some given node.

## Related nodes in a graph

Given a hyperlinked environment (= a graph)...



### Problem

Finding nodes semantically related to some given node.

# Example of related nodes

## Example (World Wide Web)

Nodes: Web pages

Edges: hyperlinks

Related nodes: similar/related pages (cf Google)

## Example (Wikipedia)

Nodes: articles

Edges: hyperlinks

Related nodes: related articles (= articles on semantically related
topics)

## Example of related nodes

### Example (World Wide Web)

Nodes: Web pages

Edges: hyperlinks

Related nodes: similar/related pages (cf Google)

### Example (Wikipedia)

Nodes: articles

Edges: hyperlinks

Related nodes: related articles (= articles on semantically related topics)

# Example of related nodes

## Example (World Wide Web)

Nodes: Web pages

Edges: hyperlinks

Related nodes: similar/related pages (cf Google)

## Example (Wikipedia)

Nodes: articles

Edges: hyperlinks

Related nodes: related articles (= articles on semantically related topics)

## Classical approaches

Classical approaches for finding related nodes (e.g. on the World Wide Web):

- Based on the use of variants of PageRank on local subgraphs.
- Text Mining techniques : cocitations, vector-space model...

### Our approach

Use of a classical Markov chain tool: Green measures.

## Classical approaches

Classical approaches for finding related nodes (e.g. on the World Wide Web):

- Based on the use of variants of PageRank on local subgraphs.
- Text Mining techniques : cocitations, vector-space model...

### Our approach

Use of a classical Markov chain tool: Green measures.

## Classical approaches

Classical approaches for finding related nodes (e.g. on the World Wide Web):

- Based on the use of variants of PageRank on local subgraphs.
- Text Mining techniques : cocitations, vector-space model...

### Our approach

Use of a classical Markov chain tool: Green measures.

## Classical approaches

Classical approaches for finding related nodes (e.g. on the World Wide Web):

- Based on the use of variants of PageRank on local subgraphs.
- Text Mining techniques : cocitations, vector-space model...

### Our approach

Use of a classical Markov chain tool: Green measures.

## Contributions

### Our contributions:

1. A novel use of Green measures for extracting semantic information in a graph.

2. An extensive comparative study with classical approaches, on the English version of Wikipedia.

### Remark

*Only pure mathematical methods, no Wikipedia-specific tricks included.*

## Contributions

Our contributions:

1. A novel use of Green measures for extracting semantic information in a graph.

2. An extensive comparative study with classical approaches, on the English version of Wikipedia.

### Remark

*Only pure mathematical methods, no Wikipedia-specific tricks included.*

## Contributions

Our contributions:

1. A novel use of Green measures for extracting semantic information in a graph.
2. An extensive comparative study with classical approaches, on the English version of Wikipedia.

### Remark

*Only pure mathematical methods, no Wikipedia-specific tricks included.*

## Contributions

Our contributions:

1. A novel use of Green measures for extracting semantic information in a graph.

2. An extensive comparative study with classical approaches, on the English version of Wikipedia.

### Remark

*Only pure mathematical methods, no Wikipedia-specific tricks included.*

# Outline

# Graph = Markov chain

### Definition (Markov chain)

Probabilistic process on a state space, defined by transition probabilities $p_{ij}$ from each state $i$ to each state $j$.

For a directed graph:

State space: set of nodes

Transition probabilities: based on existence (and weight) of edges

### Remark

All graphs will be supposed strongly connected and with gcd of length of all cycles equal to 1.

# Graph = Markov chain

### Definition (Markov chain)

Probabilistic process on a state space, defined by transition probabilities $p_{ij}$ from each state $i$ to each state $j$.

For a directed graph:

State space: set of nodes

Transition probabilities: based on existence (and weight) of edges

### Remark

All graphs will be supposed strongly connected and with gcd of length of all cycles equal to 1.

Introduction
оооо

Green measures
●оооо

Methods Compared
○○○○○○○○○○○○○○

Experiment on Wikipedia
оооо

Conclusion
оо

# Graph = Markov chain

### Definition (Markov chain)

Probabilistic process on a state space, defined by transition probabilities $p_{ij}$ from each state $i$ to each state $j$.

For a directed graph:

State space: set of nodes

Transition probabilities: based on existence (and weight) of edges

### Remark

*All graphs will be supposed strongly connected and with gcd of length of all cycles equal to 1.*

Introduction
oooo

**Green measures**
●oooo

Methods Compared
oooooooooooooo

Experiment on Wikipedia
oooo

Conclusion
oo

# Graph = Markov chain

### Definition (Markov chain)

Probabilistic process on a state space, defined by transition probabilities $p_{ij}$ from each state $i$ to each state $j$.

For a directed graph:

State space: set of nodes

Transition probabilities: based on existence (and weight) of edges

### Remark

*All graphs will be supposed strongly connected and with gcd of length of all cycles equal to 1.*

## Equilibrium measure

### Definition (Measure)

Assignments of real numbers to the state set.
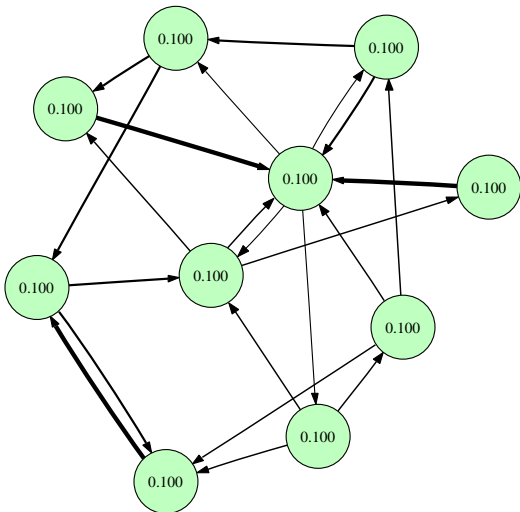
### Definition (Propagation operator)

Operator which maps a measure $\mu$ to a measure $\mu'$ computed as:

$$\mu'_j = \sum_i (\mu_i p_{ij})$$

### Result

*If we iterate the propagation operator from any measure summing to 1, we converge to a unique equilibrium measure. (PageRank with no random jumps).*

## Equilibrium measure

### Definition (Measure)

Assignments of real numbers to the state set.

### Definition (Propagation operator)

Operator which maps a measure $\mu$ to a measure $\mu'$ computed as:
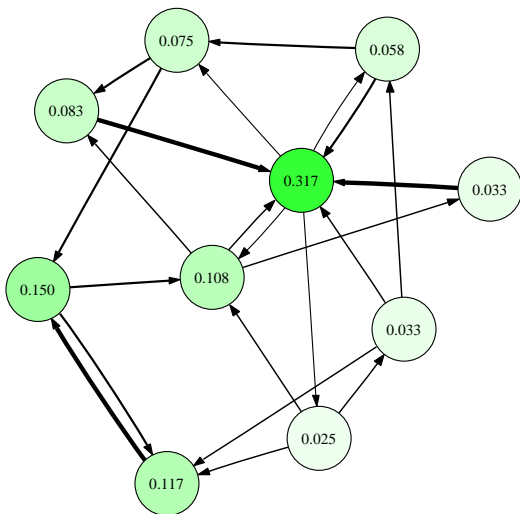
$$\mu'_j = \sum_i (\mu_i p_{ij})$$

### Result

*If we iterate the propagation operator from any measure summing to 1, we converge to a unique equilibrium measure. (PageRank with no random jumps).*
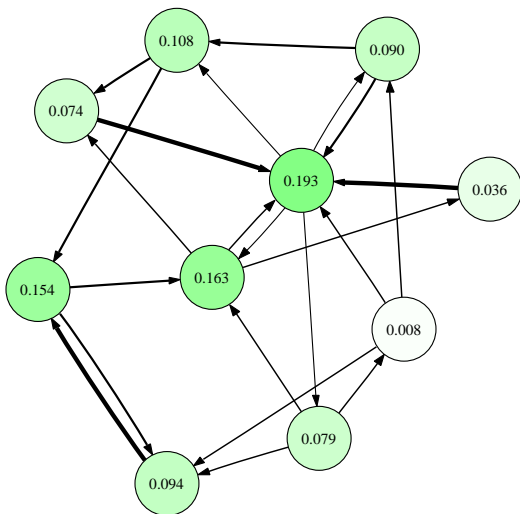
Introduction
○○○○

Green measures
○●○○○

Methods Compared
○○○○○○○○○○○○○

Experiment on Wikipedia
○○○○

Conclusion
○○

## Equilibrium measure

### Definition (Measure)

Assignments of real numbers to the state set.

### Definition (Propagation operator)

Operator which maps a measure $\mu$ to a measure $\mu'$ computed as:

$$\mu'_j = \sum_i (\mu_i p_{ij})$$

### Result

*If we iterate the propagation operator from any measure summing to 1, we converge to a unique equilibrium measure. (PageRank with no random jumps).*

Introduction
0000

**Green measures**
00●00

Methods Compared
0000000000000

Experiment on Wikipedia
0000

Conclusion
00

## $PageRank$ — Iteration ♯1

Introduction
oooo

**Green measures**
oo●oo

Methods Compared
oooooooooooooo

Experiment on Wikipedia
oooo

Conclusion
oo

## $PageRank$ — Iteration $\natural 2$

Introduction
oooo

**Green measures**
oo●oo

Methods Compared
oooooooooooooo

Experiment on Wikipedia
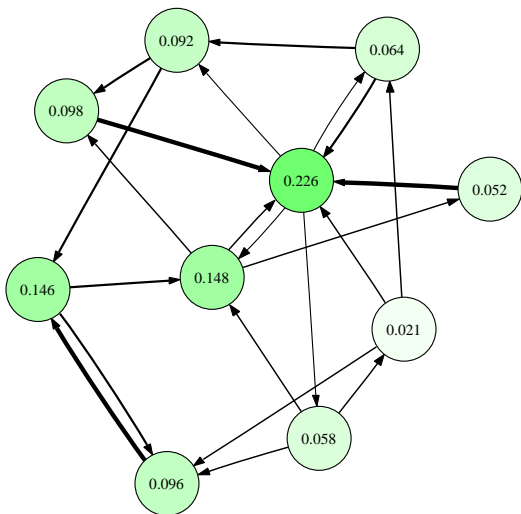oooo

Conclusion
oo

## *PageRank* — Iteration ♯3

# $PageRank$ — Iteration ♮4

Introduction
oooo

**Green measures**
oo●oo

Methods Compared
ooooooooooooo

Experiment on Wikipedia
oooo

Conclusion
oo

# *PageRank* — Iteration ♯5

Introduction
○○○○

**Green measures**
○○●○○

Methods Compared
○○○○○○○○○○○○

Experiment on Wikipedia
○○○○

Conclusion
○○

# $PageRank$ — Iteration ♮6

Introduction
oooo
**Green measures**
oo●oo
Methods Compared
oooooooooooooo
Experiment on Wikipedia
oooo
Conclusion
oo

# *PageRank* — Iteration ♯7

# *PageRank* — Iteration ♯8

Introduction
oooo

**Green measures**
ooo●o

Methods Compared
oooooooooooooo

Experiment on Wikipedia
oooo

Conclusion
oo

## $PageRank$ — Iteration ♯9

Introduction
oooo

**Green measures**
oo●oo

Methods Compared
oooooooooooooo

Experiment on Wikipedia
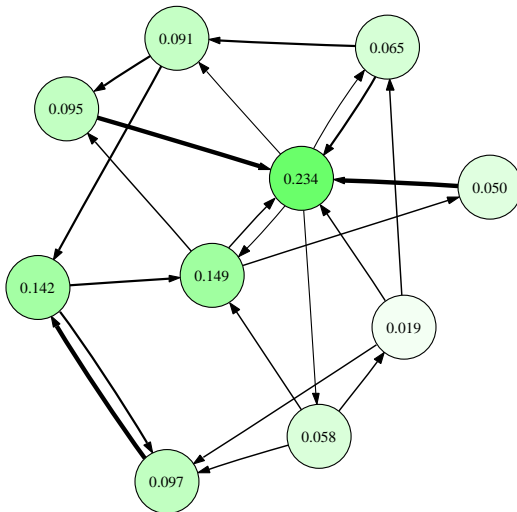oooo

Conclusion
oo

## $PageRank$ — Iteration ♮10

## $PageRank$ — Iteration $\sharp 11$

Introduction
oooo

**Green measures**
oo●oo

Methods Compared
oooooooooooooo

Experiment on Wikipedia
oooo

Conclusion
oo

## *PageRank* — Iteration ♮12

Introduction
○○○○

**Green measures**
○○●○○

Methods Compared
○○○○○○○○○○○○○

Experiment on Wikipedia
○○○○

Conclusion
○○

## $PageRank$ — Iteration ♮13

## $PageRank$ — Iteration ♮14

# Background on Green measures

## Green functions

- Come from electrostatic theory (potential created by a charge distribution).

- Analogy between electrostatic potential theory and Markov chains.

- Green measures: discrete analogues of Green functions.

## Background on Green measures

### Green functions

- Come from electrostatic theory (potential created by a charge distribution).

- Analogy between electrostatic potential theory and Markov chains.

- Green measures: discrete analogues of Green functions.

# Background on Green measures

## Green functions

- Come from electrostatic theory (potential created by a charge distribution).
- Analogy between electrostatic potential theory and Markov chains.
- Green measures: discrete analogues of Green functions.

## Definition of Green measures

### Definition (Green measure centered at node $i$)

Only <span style="color:red">fixed point</span> of the operator on measures defined by:

$$\mu_j \mapsto \sum_k (\mu_k p_{kj}) + (\delta_{ij} - \nu_j) \quad \text{where} \quad \delta_{ij} = \begin{cases} 1 \text{ if } i = j \\ 0 \text{ otherwise} \end{cases}$$

### Interpretations

- PageRank with source at $i$: standard PageRank computation while, at each iteration, adding 1 to the measure of $i$, and subtracting $\nu_j$ to every node $j$.
- Time spent at a node knowing the initial node is $i$.

Introduction
○○○○

**Green measures**
○○○○○●

Methods Compared
○○○○○○○○○○○○○

Experiment on Wikipedia
○○○○

Conclusion
○○

# Definition of Green measures

## Definition (Green measure centered at node $i$)

Only fixed point of the operator on measures defined by:

$$\mu_j \mapsto \sum_k (\mu_k p_{kj}) + (\delta_{ij} - \nu_j) \quad \text{where} \quad \delta_{ij} = \begin{cases} 1 \text{ if } i = j \\ 0 \text{ otherwise} \end{cases}$$

## Interpretations

- PageRank with source at $i$: standard PageRank computation while, at each iteration, adding 1 to the measure of $i$, and subtracting $\nu_j$ to every node $j$.

- Time spent at a node knowing the initial node is $i$.

# Definition of Green measures

## Definition (Green measure centered at node $i$)

Only <span style="color:red">fixed point</span> of the operator on measures defined by:

$$\mu_j \mapsto \sum_k (\mu_k p_{kj}) + (\delta_{ij} - \nu_j) \quad \text{where} \quad \delta_{ij} = \begin{cases} 1 \text{ if } i = j \\ 0 \text{ otherwise} \end{cases}$$

## Interpretations

- <span style="color:red">PageRank with source</span> at $i$: standard PageRank computation while, at each iteration, adding 1 to the measure of $i$, and subtracting $\nu_j$ to every node $j$.
- Time spent at a node knowing the initial node is $i$.

Introduction
oooo

**Green measures**
ooooo●

Methods Compared
oooooooooooooo

Experiment on Wikipedia
oooo

Conclusion
oo

# Definition of Green measures

## Definition (Green measure centered at node $i$)

Only <span style="color:red">fixed point</span> of the operator on measures defined by:

$$\mu_j \mapsto \sum_k (\mu_k p_{kj}) + (\delta_{ij} - \nu_j) \quad \text{where} \quad \delta_{ij} = \begin{cases} 1 \text{ if } i = j \\ 0 \text{ otherwise} \end{cases}$$

## Interpretations

- <span style="color:red">PageRank with source</span> at $i$: standard PageRank computation while, at each iteration, adding 1 to the measure of $i$, and subtracting $\nu_j$ to every node $j$.

- <span style="color:red">Time spent at a node</span> knowing the initial node is $i$.

# Outline

## Purpose

- Finding nodes in the graph related to $i$.

- For each method, output an ordered list of nodes related to $i$.

- Each method provides a similarity score to $i$.

## Purpose

- Finding nodes in the graph related to $i$.
- For each method, output an ordered list of nodes related to $i$.
- Each method provides a similarity score to $i$.

## Purpose

- Finding nodes in the graph related to $i$.
- For each method, output an ordered list of nodes related to $i$.
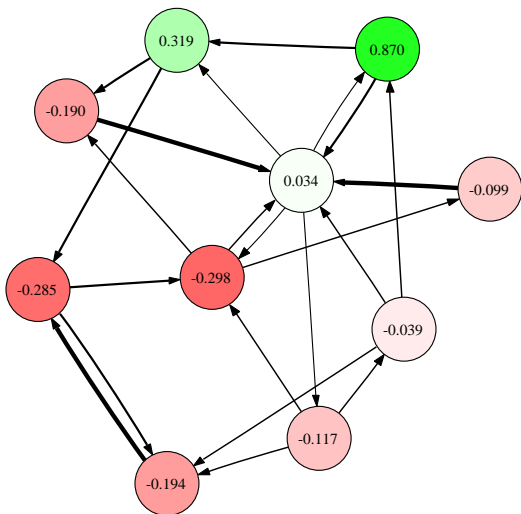- Each method provides a similarity score to $i$.

## *Green* — Method Description

### Method Description

- Direct application of the theory of Green measures.

- Improvement: multiplication by a term favoring uncommon nodes $\log(1/\nu_j)$ (quantity of information).

- Iteration until reasonable convergence on the top results.

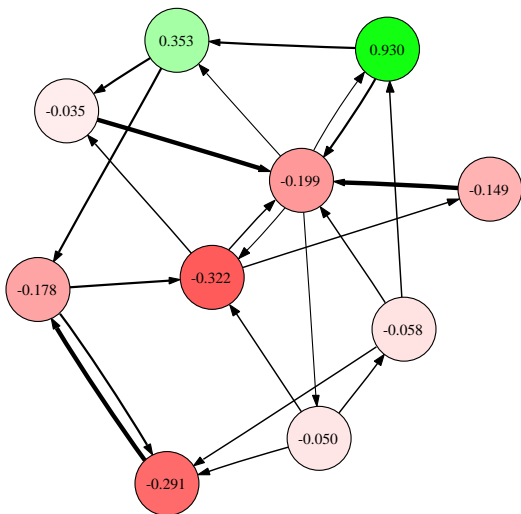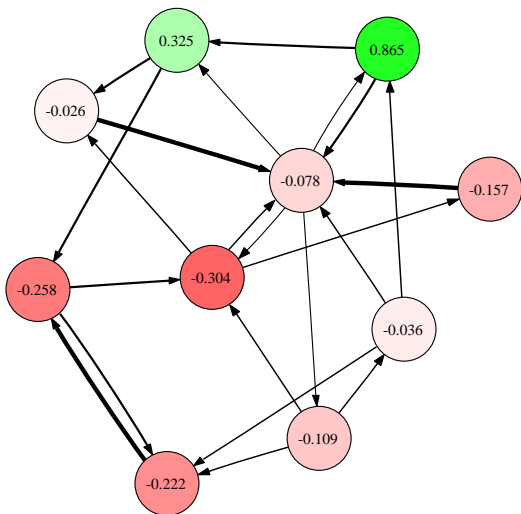## *Green* — Method Description

### Method Description

- Direct application of the theory of Green measures.
- Improvement: multiplication by a term favoring uncommon nodes $\log(1/\nu_j)$ (quantity of information).
- Iteration until reasonable convergence on the top results.
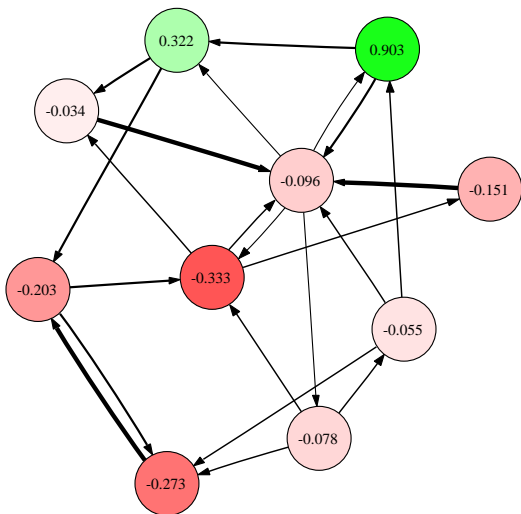
## *Green* — Method Description

### Method Description

- Direct application of the theory of Green measures.
- Improvement: multiplication by a term favoring uncommon nodes $\log(1/\nu_j)$ (quantity of information).
- Iteration until reasonable convergence on the top results.

Introduction
0000

Green measures
00000

**Methods Compared**
0●00000000000

Experiment on Wikipedia
0000

Conclusion
00

## *Green* — Iteration ♮1

Introduction
0000

Green measures
00000

Methods Compared
0●0000000000

Experiment on Wikipedia
0000

Conclusion
00

## *Green* — Iteration ♮2

Introduction
OOOO

Green measures
OOOOO

Methods Compared
O●OOOOOOOOOOO

Experiment on Wikipedia
OOOO

Conclusion
OO

## *Green* — Iteration ♮3

## *Green* — Iteration ♮4

Introduction
oooo

Green measures
ooooo

**Methods Compared**
o●oooooooooooo

Experiment on Wikipedia
oooo

Conclusion
oo

## *Green* — Iteration ♮5

## *Green* — Iteration ♮6

Introduction
oooo

Green measures
ooooo

**Methods Compared**
○●○○○○○○○○○○○

Experiment on Wikipedia
oooo

Conclusion
oo

## *Green* — Iteration ♯7

Introduction
oooo

Green measures
ooooo

Methods Compared
o●ooooooooooo

Experiment on Wikipedia
oooo

Conclusion
oo

## *Green* — Iteration ♯8

Introduction
○○○○

Green measures
○○○○○

**Methods Compared**
○●○○○○○○○○○○○

Experiment on Wikipedia
○○○○

Conclusion
○○

# *Green* — Iteration ♯9

Introduction
oooo
Green measures
ooooo
Methods Compared
o●oooooooooo
Experiment on Wikipedia
oooo
Conclusion
oo

## *Green* — Iteration ♮10

Introduction
0000

Green measures
00000

**Methods Compared**
0●00000000000

Experiment on Wikipedia
0000

Conclusion
00

## *Green* — Iteration ♮11

Introduction
0000

Green measures
00000

**Methods Compared**
0●0000000000

Experiment on Wikipedia
0000

Conclusion
00

## *Green* — Iteration ♮12

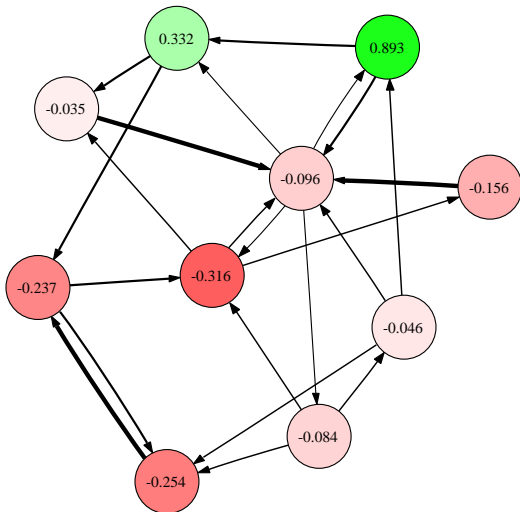## *Green* — Iteration ♮13

## *Green* — Iteration ♮14

Introduction
0000

Green measures
00000

Methods Compared
0●00000000000

Experiment on Wikipedia
0000

Conclusion
00

## *Green* — Iteration ♮15

## *Green* — Iteration ♮16

Introduction
0000

Green measures
00000

**Methods Compared**
○●○○○○○○○○○○○

Experiment on Wikipedia
0000

Conclusion
00

## *Green* — Iteration ♮17

Introduction
0000
Green measures
00000
**Methods Compared**
○●○○○○○○○○○○○
Experiment on Wikipedia
0000
Conclusion
00

## *Green* — Iteration ♮18

Introduction
oooo

Green measures
ooooo

**Methods Compared**
○●○○○○○○○○○○○

Experiment on Wikipedia
oooo

Conclusion
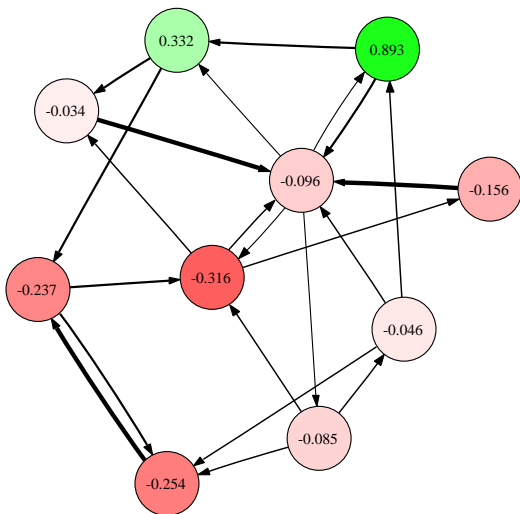oo

## *Green* — Iteration ♮19

## *Green* — Iteration ♮20

# *SymGreen* — Method Description

## Method Description

- *Green* goes only forward, may be a limitation.
- Symmetrize the graph, in a canonical sense in relation to the equilibrium measure:

$$\tilde{p}_{ij} = \frac{1}{2}(p_{ij} + p_{ji}\frac{\nu_j}{\nu_i})$$

  The resulting graph has the same equilibrium measure.

- Same as *Green* on this symmetrized graph.

Introduction
oooo

Green measures
ooooo

**Methods Compared**
oo●oooooooooooo

Experiment on Wikipedia
oooo

Conclusion
oo

## *SymGreen* — Method Description

### Method Description

- *Green* goes only forward, may be a limitation.
- Symmetrize the graph, in a canonical sense in relation to the equilibrium measure:

$$\tilde{p}_{ij} = \frac{1}{2}\left(p_{ij} + p_{ji}\frac{\nu_j}{\nu_i}\right)$$

  The resulting graph has the same equilibrium measure.

- Same as *Green* on this symmetrized graph.
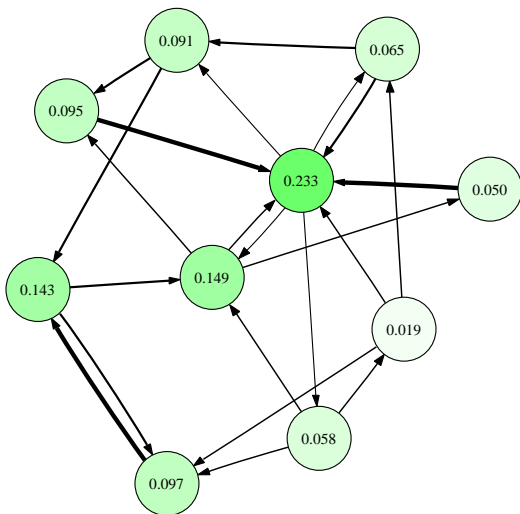
## *SymGreen* — Method Description

### Method Description

- *Green* goes only forward, may be a limitation.
- Symmetrize the graph, in a canonical sense in relation to the equilibrium measure:

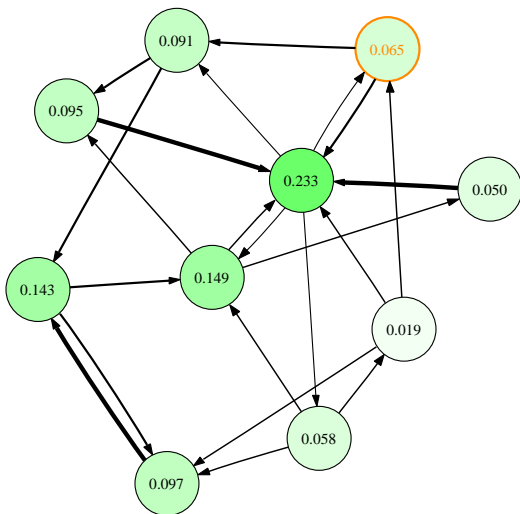$$\tilde{p}_{ij} = \frac{1}{2}\left(p_{ij} + p_{ji}\frac{\nu_j}{\nu_i}\right)$$

  The resulting graph has the same equilibrium measure.
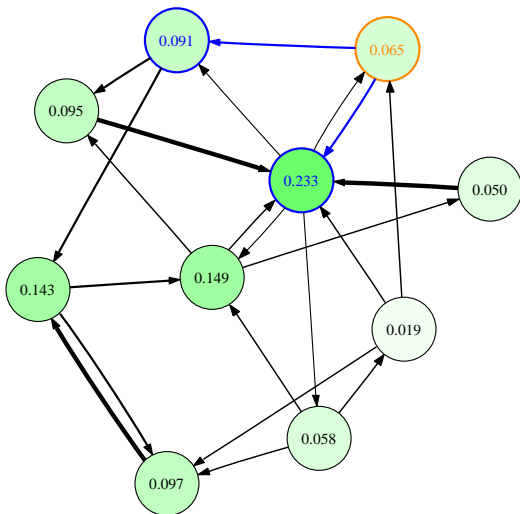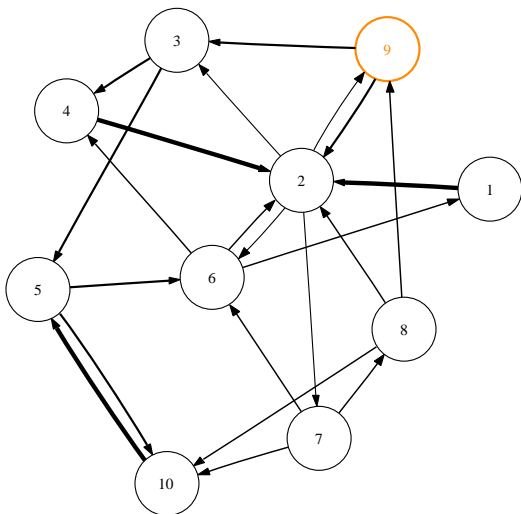- Same as *Green* on this symmetrized graph.

## PageRankOfLinks

## PageRankOfLinks

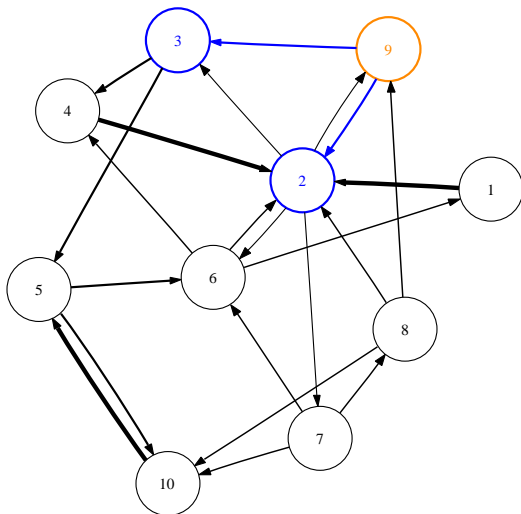Introduction
oooo

Green measures
ooooo

**Methods Compared**
oooooo●oooooo

Experiment on Wikipedia
oooo

Conclusion
oo

## PageRankOfLinks

## Cosine

# Cosine

Introduction
0000

Green measures
00000

Methods Compared
000000000●0000

Experiment on Wikipedia
0000

Conclusion
00

## *Cosine*

## Cosine

Dimensions

| | 1 | 2 | 3 | 4 | 6 | 7 | 9 | 10 | Cosine with 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | | ✓ | | | | | | | 0.40 |
| 2 | | | ✓ | | ✓ | ✓ | ✓ | | 0.43 |
| 4 | | ✓ | | | | | | | 0.40 |
| 6 | ✓ | ✓ | | ✓ | | | | | 0.09 |
| 8 | | ✓ | | | | | ✓ | ✓ | 0.13 |
| 9 | | ✓ | ✓ | | | | | | 1.00 |

Documents

Introduction
0000

Green measures
00000

**Methods Compared**
0000000000●00

Experiment on Wikipedia
0000

Conclusion
00

## Cocitations

## Cocitations

## Cocitations

# Outline

## The graph of Wikipedia

### Statistics

- $1, 606, 896$ nodes (as of September 25th, 2006).

- $38, 896, 462$ edges.

- $95\%$ of the nodes belong to the largest strongly connected component.

# Evaluation methodology

- Blind evaluation of the methods.

- Articles selected for their diversity:

    - Clique (graph theory)

    - Germany

    - Hungarian language

    - Pierre de Fermat

    - Star Wars

    - Theory of relativity

    - 1989

- 66 evaluators asked to give a mark to each list of words.

## Evaluation methodology

- Blind evaluation of the methods.
- Articles selected for their diversity:
  - Clique (graph theory)
  - Germany
  - Hungarian language
  - Pierre de Fermat
  - Star Wars
  - Theory of relativity
  - 1989
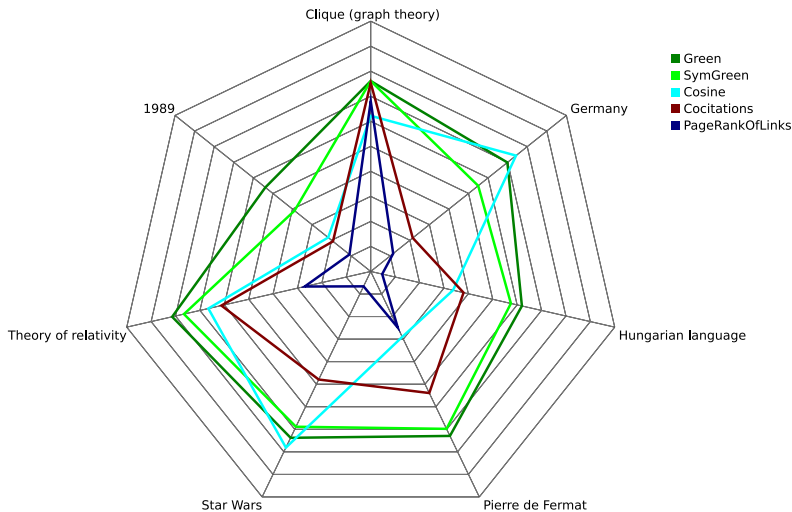- 66 evaluators asked to give a mark to each list of words.

## Evaluation methodology

- Blind evaluation of the methods.
- Articles selected for their diversity:
  - Clique (graph theory)
  - Germany
  - Hungarian language
  - Pierre de Fermat
  - Star Wars
  - Theory of relativity
  - 1989
- 66 evaluators asked to give a mark to each list of words.

## Output on Germany

| Green | SymGreen | PageRankOfLinks | Cosine | Cocitations |
|---|---|---|---|---|
| 1. Germany<br>2. Berlin<br>3. German language<br>4. Christian Democratic Union (Germany)<br>5. Austria<br>6. Hamburg<br>7. German reunification<br>8. Social Democratic Party of Germany<br>9. German Empire<br>10. German Democratic Republic | 1. Germany<br>2. Berlin<br>3. France<br>4. Austria<br>5. German language<br>6. Bavaria<br>7. World War II<br>8. German Democratic Republic<br>9. European Union<br>10. Hamburg | 1. United States<br>2. United Kingdom<br>3. France<br>4. 2005<br>5. Germany<br>6. World War II<br>7. Canada<br>8. English language<br>9. Japan<br>10. Italy | 1. Germany<br>2. History of Germany since 1945<br>3. History of Germany<br>4. Timeline of German history<br>5. States of Germany<br>6. Politics of Germany<br>7. List of Germany-related topics<br>8. Hildesheimer Rabbinical Seminary<br>9. Pleasure Victim<br>10. German Unity Day | 1. Germany<br>2. United States<br>3. France<br>4. United Kingdom<br>5. World War II<br>6. Italy<br>7. Netherlands<br>8. Japan<br>9. 2005<br>10. Category:Living people |

Introduction
oooo

Green measures
ooooo

Methods Compared
oooooooooooooo

Experiment on Wikipedia
ooo●

Conclusion
oo

## Results

# Outline

## Summary

- Green measures: a tool for extracting semantic information in a graph.
- In comparison to other methods, in the case of Wikipedia:
  - Better overall performance.
  - Robustness.
  - Discovery of relevant semantic relations.

| Introduction | Green measures | Methods Compared | Experiment on Wikipedia | Conclusion |
|:---|:---|:---|:---|:---|
| oooo | ooooo | oooooooooooo | oooo | ●o |

## Summary

- Green measures: a tool for extracting semantic information in a graph.
- In comparison to other methods, in the case of Wikipedia:
  - Better overall performance.
  - Robustness.
  - Discovery of relevant semantic relations.

## Summary

- Green measures: a tool for extracting semantic information in a graph.
- In comparison to other methods, in the case of Wikipedia:
  - Better overall performance.
  - Robustness.
  - Discovery of relevant semantic relations.

## Summary

- Green measures: a tool for extracting semantic information in a graph.
- In comparison to other methods, in the case of Wikipedia:
  - Better overall performance.
  - Robustness.
  - Discovery of relevant semantic relations.

Introduction
oooo

Green measures
ooooo

Methods Compared
oooooooooooooo

Experiment on Wikipedia
oooo

Conclusion
●o

## Summary

- Green measures: a tool for extracting semantic information in a graph.
- In comparison to other methods, in the case of Wikipedia:
  - Better overall performance.
  - Robustness.
  - Discovery of relevant semantic relations.

# Perspectives



- Application to the Web graph.

- Interpolation between *Green* and *SymGreen*.

- Clustering using *Green* measures: unpractical now because of computation times.

- Use of *Green* measures on other Markov chains, e.g. for computing authority scores.

# Perspectives



- Application to the Web graph.

- Interpolation between *Green* and *SymGreen*.

- Clustering using *Green* measures: unpractical now because of computation times.

- Use of *Green* measures on other Markov chains, e.g. for computing authority scores.

## Perspectives



- Application to the Web graph.

- Interpolation between *Green* and *SymGreen*.

- Clustering using *Green* measures: unpractical now because of computation times.

- Use of *Green* measures on other Markov chains, e.g. for computing authority scores.

## Perspectives



- Application to the Web graph.
- Interpolation between *Green* and *SymGreen*.
- Clustering using *Green* measures: unpractical now because of computation times.
- Use of *Green* measures on other Markov chains, e.g. for computing authority scores.