# Pagerank in Apache Flink

Author: Ward Schodts
Supervisor: Juan Soto

Datenbanksysteme und Informationsmanagement
Technische Universität Berlin

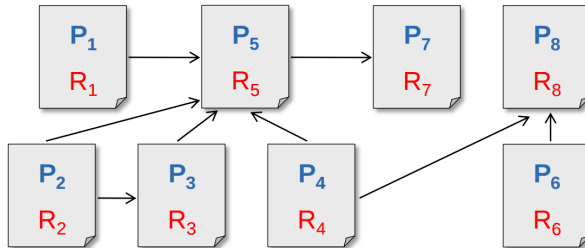July 8, 2016

## Agenda

Introduction

The experiment

The different algorithm implementations

Results

Conclusion

# Pagerank



PageRank example 1 [5]

- A page has a high PageRank *R* if
  - there are many pages linking to it
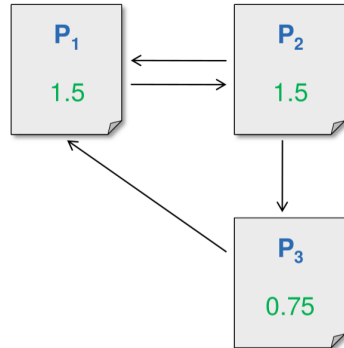  - or, if there are some pages with a high PageRank linking to it

## Pagerank

$$R(P_i) = \sum_{P_j \in B_i} \frac{R(P_j)}{L_j}$$

▶ where

- $B_i$ is the set of pages that link to page $P_i$
- $L_j$ is the number of outgoing links for page $P_j$ linking to it



PageRank example 2 [5]

[2]

## Apache Flink

- ▶ Open source framework for distributed Big Data Analytics
- ▶ Exploits:
    - data streaming
    - in-memory processing
    - iteration operators

    to improve performance
- ▶ Formerly Stratosphere (Flink means agile)
- ▶ Developped here at TUB

[6]

## Apache Flink: 2 possible setups

```xml
<dependencies>

    <dependency>
        <groupId>org.apache.flink</groupId>
        <artifactId>flink-java</artifactId>
        <version>${flink.version}</version>
    </dependency>


    <dependency>
        <groupId>org.apache.flink</groupId>
        <artifactId>flink-streaming-java_2.10</artifactId>
        <version>${flink.version}</version>
    </dependency>


    <dependency>
        <groupId>org.apache.flink</groupId>
        <artifactId>flink-gelly_2.10</artifactId>
        <version>${flink.version}</version>
    </dependency>

    <dependency>
        <groupId>org.apache.flink</groupId>
        <artifactId>flink-table_2.10</artifactId>
        <version>${flink.version}</version>
    </dependency>
```



Binary version (self compiled)

Maven

## Agenda

## General experiment setup

1. Data file with graph and pagerank solution

# Flink experiment setup

# GraphLab experiment setup

## Agenda

# Flink algorithm 1

# Flink algorithm 2

# Flink algorithm 3

# GraphLab algorithm

# Agenda

## Results

## Agenda

## Conclusion

## Thank you for your attention

Questions?

References I

📄 Slim Baltagi. *Overview of Apache Flink: Next-Gen Big Data Analytics Framework*. 2015. URL: %5Curl%7Bhttp: //www.slideshare.net/sbaltagi/overview-of-apacheflinkbyslimbaltagi?qid=5f0b5424-d187-4c79-a600-6cae794c686e&v=&b=&from_search=3%7D.

📄 Apache Flink. *Apache Flink Squirrel*. URL: %5Curl%7Bhttps://flink.apache.org/img/logo/png/1000/flink_squirrel_1000.png%7D.

📄 Lawrence Page et al. "The PageRank citation ranking: bringing order to the web." In: (1999).

📄 Lawrence Page et al. "The PageRank citation ranking: bringing order to the web." In: (1999).

References II

📄 Beat signer. *Google PageRank*. 2009. URL: %5Curl%7Bhttp:
//www.slideshare.net/signer/google-
pagerank-presentation?qid=18af8836-30e7-
41cd-9edb-
956bd7ca324d&v=&b=&from_search=2%7D.

📄 Mathias Spahlinger. *There is no repetition.* URL:
%5Curl%7Bhttps://www.google.com/search?q=
repeat&source=lnms&tbm=isch&sa=X&ved=
0ahUKEwi4laH2tuLNAhVnB8AKHTPQCU4Q_AUICCgB&
biw=1590&bih=765#tbm=isch&q=no+repetition&
imgrc=h1qwLbEEezv8SM:%7D.