

# Text Technologies for Data Science Assignment 4

S1140740

November 13, 2014

## Introduction

### PageRank Algorithm

People are represented as nodes with a directed edge from a person sending an email to a person receiving the email. The weight ( $w$ ) of this edge represents the number of emails the person sent. PageRanks ( $PR$ ) are initialised to  $1/N$  where  $N$  is the total number of people. The first 10 iterations of the algorithm are then executed. Each iteration  $S$ , the sum of PageRanks of sink nodes, is computed and people's PageRanks are updated using the PageRank formula.<sup>1</sup>

Note that every iteration PageRanks are computed from PageRanks from the previous iteration.

<sup>1</sup>

$$PR(x) = \frac{1 - \lambda + \lambda S}{N} + \lambda \sum_{y \rightarrow x} \frac{w \cdot PR(y)}{out(y)}$$

### HITS Algorithm

Same representation as for PageRank is used. However, hub and authority values are initialised to  $1/\sqrt{N}$ . Again, first 10 iterations of the algorithm are executed. Each iteration the hub value of a node is updated using authority values of nodes it is pointing to it<sup>2</sup> and the authority value is updated using hub values pointing to it.<sup>3</sup>

Both hub and authority values are then normalized<sup>4</sup>.

<sup>2</sup>

$$H(x) = \sum_{y \leftarrow x} A(y)$$

<sup>3</sup>

$$A(x) = \sum_{y \rightarrow x} H(y)$$

<sup>4</sup>

$$\sum_x H(x)^2 = 1 = \sum_x A(x)^2$$