

## Stochastische Prozesse

### Woche 5

#### Aufgabe 1 PageRank-Algorithmus von Google

Das Herz der Google-Suchmaschine ist ein Algorithmus, der alle Dokumente des WWW nach ihrer Wichtigkeit anordnet. Die Auflistung der Dokumente, die unter einem bestimmten Suchbegriff gefunden werden, legt diese Rangfolge, den sogenannten *PageRank*, zugrunde.<sup>1</sup>

Grundidee ist es, das Surf-Verhalten eines repräsentativen Internet-Surfers als geeignete Markov-Kette zu modellieren, wobei jede Site einem Zustand entspricht. Die stationäre Verteilung gibt dann an, wie häufig dieser Surfer (wenn er sehr lange im Netz bleibt) an den verschiedenen Sites vorbeikommt. Die Seite mit der grössten asymptotischen Wahrscheinlichkeit wird als die wichtigste interpretiert und erhält den höchsten PageRank.

Die Übergangsmatrix für ein Universum aus  $N$  Websites wird nach folgenden Regeln aufgebaut:

- Die Übergangsmatrix  $P$  ergibt sich aus zwei Matrizen  $U$  und  $S$  durch Addition:

$$P = \alpha \cdot S + (1 - \alpha) \cdot U \quad 0 \leq \alpha \leq 1$$

$P$  heisst auch *Google-Matrix*. Die Erfinder von PageRank, Sergey Brin und Larry Page, verwendeten  $\alpha = 0.85$ . Die Matrix  $S$  entspricht dem Surf-Verhalten aufgrund der Links auf einer Seite:

- Für eine Seite  $i$ , die  $l_i$  Links auf andere Seiten enthält, setzen wir die Übergangswahrscheinlichkeit

$$S_{ij} = \frac{1}{l_i},$$

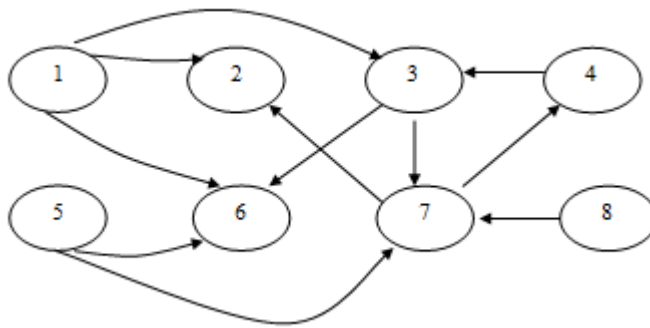
falls ein Link nach  $j$  existiert. Wenn kein Link existiert setzen wir  $S_{ij} = 0$ . Das heisst, wir nehmen an, dass ein Surfer zufällig einem der Links folgt.

- Für eine Seite, die keine Links hat, setzen wir  $S_{ij} = 1/N$  für alle  $j = 1, \dots, N$ . Das heisst, wir nehmen an, dass der Surfer zufällig irgendeine andere WebSite aufsucht.
- Die Matrix  $U$  entspricht einem rein zufälligen Surf-Verhalten ohne Links:  $U_{ij} = 1/N$  für alle  $i, j = 1 \dots, N$ .

Wenn wir das Surf-Verhalten nur aufgrund der Links einer Seite modellieren wollten, würden wir  $\alpha = 1$  verwenden. Dies entspricht der Topologie des WWW. Die Wahl  $\alpha = 0$  entspricht einem absolut zufälligen Surfen. Setzen wir  $\alpha = 0.85$ , ist das Verhalten dominiert von den weiterführenden Links, enthält aber auch andere Anteile, die wir als zufällig beschreiben.

<sup>1</sup>Literatur: David Austin, How Google finds Your Needle in the Web's Haystack. <http://www.ams.org/featurecolumn/archive/pagerank.html>

- a) Zeigen Sie allgemein, dass  $P$  für beliebiges  $\alpha$  eine zulässige Übergangsmatrix ist, d.h. dass  $0 \leq P_{ij} \leq 1$  ist, und dass die Zeilensummen gleich 1 sind.
- b) Zeigen Sie, dass  $P$  für  $\alpha < 1$  irreduzibel und aperiodisch ist. Das heisst, dass es genau eine asymptotische Verteilung gibt.
- c) Berechnen Sie die PageRanks für alle Seiten des folgenden Systems (Pfeile bedeuten Links) für  $\alpha = 1$ ,  $\alpha = 0.85$ , und  $\alpha = 0$ . Sortieren Sie die Seiten nach ihrem Rang. Ändert sich die Reihenfolge der Seiten, wenn Sie von  $\alpha = 1$  auf  $\alpha = 0.85$  gehen?



- d) Nehmen Sie an, Sie sind Eigentümer der Site 5. Sie wollen Ihr Ranking verbessern, indem Sie einen Link von einer der anderen Sites auf Ihre Site erhalten. Von welcher Site aus wäre der Link am wirkungsvollsten, d.h. bei welchem zusätzlichen Link  $i \rightarrow 5$  würden Sie Ihr Ranking am meisten verbessern? Wie wäre Ihr neuer Platz im Ranking, wenn Sie diesen Link erhalten würden? (Verwenden Sie hier  $\alpha = 0.85$ ).

## Aufgabe 2 Spezielle Zustände

Eine Markov-Kette habe die Übergangsmatrix:

$$P = \begin{pmatrix} 0.00 & 1.0 \\ 1.00 & 0.0 \end{pmatrix}$$

- a) Zeichnen Sie das Übergangsdiagramm.
- b) Plotten Sie  $\vec{\pi}(t)$  jeweils für die Anfangsverteilung  $\vec{\pi}(0) = (1, 0)$  und  $\vec{\pi}(0) = (0, 1)$  gegen die Zeit  $t$ .
- c) Ersetzen Sie nun die Übergangsmatrix durch

$$P = \begin{pmatrix} 0.05 & 0.95 \\ 0.95 & 0.05 \end{pmatrix}$$

und erstellen Sie den gleichen Plot wie in b).