



## **Spectral properties of Google matrix**

#### **Klaus Frahm**

#### **Quantware MIPS Center**

Université Paul Sabatier Laboratoire de Physique Théorique, UMR 5152, IRSAMC

#### A. D. Chepelianskii, Y. H. Eom, L. Ermann, B. Georgeot, D. Shepelyansky

Quantum chaos: fundamentals and applications Luchon, March 14 - 21, 2015

# Contents

Perron-Frobenius operators	3
PF Operators for directed networks	4
PageRank	6
Numerical diagonalization	7
University Networks	9
Wikipedia	12
Twitter network	14
Random Perron-Frobenius matrices	16
Poisson statistics of PageRank	18
Physical Review network	20
Perron-Frobenius matrix for chaotic maps	26
References	35

## **Perron-Frobenius operators**

Consider a physical system with N states i = 1, ..., N and probabilities  $p_i(t) \ge 0$  evolving by a discrete *Markov process*:

$$p_i(t+1) = \sum_j G_{ij} p_j(t)$$
 with  $\sum_i G_{ij} = 1$  ,  $G_{ij} \ge 0$ .

The transition probabilities  $G_{ij}$  provide a *Perron-Frobenius* matrix. Conservation of probability:  $\sum_i p_i(t+1) = \sum_i p_i(t) = 1$ .

In general  $G^T \neq G$  and eigenvalues  $\lambda$  may be complex and obey  $|\lambda| \leq 1$ . The vector  $e^T = (1, \ldots, 1)$  is left eigenvector with  $\lambda_1 = 1$  $\Rightarrow$  existence of (at least) one right eigenvector P for  $\lambda_1 = 1$  also called **PageRank** in the context of Google matrices: GP = 1P

For non-degenerate  $\lambda_1$  and finite gap  $|\lambda_2| < 1$ :  $\lim_{t \to \infty} p(t) = P$ 

 $\Rightarrow$  **Power method** to compute P with rate of convergence  $\sim |\lambda_2|^t$ .

# **PF Operators for directed networks**

Consider a directed network with N nodes  $1,\,\ldots,\,N$  and  $N_\ell$  links. Adjacency matrix:

 $A_{jk} = 1$  if there is a link  $k \rightarrow j$  and  $A_{jk} = 0$  otherwise.

Sum-normalization of each non-zero column of  $A \Rightarrow S_0$ .

Replacing each zero column (*dangling nodes*) with  $e/N \Rightarrow S$ .

Eventually apply the *damping factor*  $\alpha < 1$  (typically  $\alpha = 0.85$ ):

#### Google matrix:

$$G(\alpha) = \alpha S + (1 - \alpha) \frac{1}{N} e e^T \quad .$$

 $\Rightarrow \lambda_1$  is non-degenerate and  $|\lambda_2| \leq \alpha$ .

Same procedure for inverted network:  $A^* \equiv A^T$  where  $S^*$  and  $G^*$  are obtained in the same way from  $A^*$ . Note: in general:  $S^* \neq S^T$ . Leading (right) eigenvector of  $S^*$  or  $G^*$  is called *CheiRank*.

## Example:

$$S_{0} = \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{3} & 0 & 0 \\ 1 & 0 & \frac{1}{3} & \frac{1}{3} & 0 \\ 0 & \frac{1}{2} & 0 & \frac{1}{3} & \frac{1}{3} & 0 \\ 0 & \frac{1}{2} & 0 & \frac{1}{3} & 0 \\ 0 & 0 & \frac{1}{3} & 0 & 0 \\ 0 & 0 & \frac{1}{3} & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{3} & 0 \end{pmatrix} , \quad S = \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{3} & 0 & \frac{1}{5} \\ 1 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{5} \\ 0 & \frac{1}{2} & 0 & \frac{1}{3} & \frac{1}{5} \\ 0 & 0 & \frac{1}{3} & 0 & \frac{1}{5} \\ 0 & 0 & 0 & \frac{1}{3} & 0 & \frac{1}{5} \\ 0 & 0 & 0 & \frac{1}{3} & 0 & \frac{1}{5} \\ 0 & 0 & 0 & \frac{1}{3} & \frac{1}{5} \end{pmatrix}$$

# PageRank

Example for university networks of Cambridge 2006 and Oxford 2006 ( $N \approx 2 \times 10^5$  and  $N_\ell \approx 2 \times 10^6$ ).



P(i) represents the "importance" of "node/page i" obtained as sum of all other pages j pointing to i with weight P(j). Sorting of  $P(i) \Rightarrow$  index K(i) for order of appearance of search results in search engines such as Google.

## **Numerical diagonalization**

- **Power method** to obtain P: rate of convergence for  $G(\alpha) \sim \alpha^t$ .
- Full "exact" diagonalization ( $N \lesssim 10^4$ ).
- Arnoldi method to determine largest  $n_A \sim 10^2 10^4$  eigenvalues. Idea: write

$$G \xi_k = \sum_{j=0}^{k+1} H_{jk} \xi_j$$
 for  $k = 0, \dots, n_A - 1$ 

where  $\xi_{k+1}$  is obtained from *Gram-Schmidt* orthogonalization of  $G\xi_k$  to  $\xi_0, \ldots, \xi_k$  with  $\xi_0$  being some suitable normalized initial vector.  $\xi_0, \ldots, \xi_{n_A-1}$  span a *Krylov space* of dimension  $n_A$  and the eigenvalues of the "small" representation matrix  $H_{jk}$  are (very) good approximations to the largest eigenvalues of G. Example for Twitter network of 2009:  $N \approx 4 \times 10^7$  and  $N_\ell \approx 1.5 \times 10^9$  with  $n_A = 640$  (lower N in other examples allows for higher  $n_A$ ). • Practical problems due to *invariant subspaces* of nodes in realistic WWW networks creating large degeneracies of  $\lambda_1$  (or  $\lambda_2$  if  $\alpha < 1$ ). Decomposition in subspaces and a core space

$$\Rightarrow \quad S = \left(\begin{array}{cc} S_{ss} & S_{sc} \\ 0 & S_{cc} \end{array}\right)$$

where  $S_{ss}$  is block diagonal according to the subspaces. The subspace blocks of  $S_{ss}$  are all matrices of PF type with at least one eigenvalue  $\lambda_1 = 1$  explaining the high degeneracies.

To determine the spectrum of S apply exact (or Arnoldi) diagonalization on each subspace and the Arnoldi method to  $S_{cc}$  to determine the largest core space eigenvalues  $\lambda_j$  (note:  $|\lambda_j| < 1$ ).

 Strange numerical problems to determine accurately "small" eigenvalues, in particular for (nearly) *triangular network structure* due to large Jordan-blocks (e.g. citation network of Physical Review).

# **University Networks**



Cambridge 2006 (left),  $N = 212710, N_s = 48239$ Oxford 2006 (right),  $N = 200823, N_s = 30579$ 

Spectrum of S (upper panels),  $S^*$  (middle panels) and dependence of rescaled level number on  $|\lambda_j|$  (lower panels).

Blue: subspace eigenvalues Red: core space eigenvalues (with Arnoldi dimension  $n_A = 20000$ ) PageRank for  $\alpha \rightarrow 1$  :



10

## Core space gap and quasi-subspaces



Left: Core space gap  $1 - \lambda_1^{(\text{core})}$  vs N for certain british universities. Red dots for gap  $> 10^{-9}$ ; blue crosses (moved up by  $10^9$ ) for gap  $< 10^{-16}$ . Right: first core space eigenvecteur for universities with gap  $< 10^{-16}$  or gap  $= 2.91 \times 10^{-9}$  for Cambridge 2004.

Core space gaps  $< 10^{-16}$  correspond to *quasi-subspaces* where it takes quite many "iterations" to reach a dangling node.

# Wikipedia

Wikipedia 2009 : N = 3282257 nodes,  $N_{\ell} = 71012307$  network links.



left (right): PageRank (CheiRank)

<u>black:</u> PageRank (CheiRank) at  $\alpha = 0.85$ <u>grey:</u> PageRank (CheiRank) at  $\alpha = 1 - 10^{-8}$ red and green: first two core space eigenvectors

blue and pink: two eigenvectors with large imaginary part in the eigenvalue



#### "Themes" of certain Wikipedia eigenvectors:

## **Twitter network**

Twitter 2009 : N=41652230 nodes,  $N_\ell=1468365182$  network links.

Matrix structure in K-rank order:



Number  $N_G$  of non-empty matrix elements in  $K \times K$ -square:



#### **Spectrum for the Twitter network**



 $n_A = 640 \implies$  requires  $\sim 200$  GB of RAM memory.

# Random Perron-Frobenius matrices

Construct random matrix ensembles  $G_{ij}$  such that:

 $G_{ij} \ge 0$ ,  $G_{ij}$  are (approximately) non-correlated and distributed with the same distribution  $P(G_{ij})$  (of finite variance  $\sigma^2$ ),

$$\sum_{j} G_{ij} = 1 \quad \Rightarrow \quad \langle G_{ij} \rangle = 1/N$$

 $\Rightarrow$  average of *G* has one eigenvalue  $\lambda_1 = 1$  ( $\Rightarrow$  "flat" PageRank) and other eigenvalues  $\lambda_j = 0$  (for  $j \neq 1$ ).

degenerate perturbation theory for the fluctuations  $\Rightarrow$  circular eigenvalue density with  $R = \sqrt{N}\sigma$  and one unit eigenvalue.

Different variants of the model:

full  $\Rightarrow$   $R = 1/\sqrt{3N}$ 

**sparse** with Q non-zero elements per column  $\Rightarrow R \sim 1/\sqrt{Q}$ 

power law with  $P(G) \sim G^{-b}$  for  $2 < b < 3 \implies R \sim N^{1-b/2}$ 

#### Numerical verification:



# **Poisson statistics of PageRank**



Identify PageRank values to "energy-levels":

$$P(i) = \exp(-E_i/T)/Z$$

with  $Z = \sum_{i} \exp(-E_i/T)$  and an effective temperature T (can be choosen: T = 1).



Parameter dependance of  $E_i = -\ln(P(i))$  on the damping factor  $\alpha$ .

# **Physical Review network**

 $N=463347 \ \mathrm{nodes}$  and  $N_\ell=4691015 \ \mathrm{links}.$ 

Coarse-grained matrix structure ( $500 \times 500$  cells):



left: time ordered, right: journal and then time ordered

"11" Journals of Physical Review: (Phys. Rev. Series I), Phys. Rev., Phys. Rev. Lett., (Rev. Mod. Phys.), Phys. Rev. A, B, C, D, E, (Phys. Rev. STAB and Phys. Rev. STPER).

 $\Rightarrow$  nearly triangular matrix structure of adjacency matrix: most citations links  $t \to t'$  are for t > t' ("past citations") but there is a small number ( $12126 = 2.6 \times 10^{-3} N_{\ell}$ ) of links  $t \to t'$  with  $t \le t'$  corresponding to *future citations*.

Strong numerical problems due to large Jordan subspaces!

## **Triangular approximation**

Remove the small number of links due to "future citations".

Semi-analytical diagonalization is possible:

$$S = S_0 + e \, d^T / N$$

where  $e_n = 1$  for all nodes n,  $d_n = 1$  for dangling nodes n and  $d_n = 0$  otherwise.  $S_0$  is the pure link matrix which is *nil-potent*:

$$S_0^l = 0$$
 with  $l = 352$ .

Let  $\psi$  be an eigenvector of S with eigenvalue  $\lambda$  and  $C = d^T \psi$ . If  $C = 0 \Rightarrow \psi$  eigenvector of  $S_0 \Rightarrow \lambda = 0$  since  $S_0$  nil-potent.

These eigenvectors belong to large Jordan blocks and are responsible for the numerical problems.

If  $C \neq 0 \Rightarrow \lambda \neq 0$  since the equation  $S_0\psi = -C e/N$  does not have a solution  $\Rightarrow \lambda \mathbf{1} - S_0$  invertible.

$$\Rightarrow \psi = C \left(\lambda \mathbf{1} - S_0\right)^{-1} e/N = \frac{C}{\lambda} \sum_{j=0}^{l-1} \left(\frac{S_0}{\lambda}\right)^j e/N$$
  
From  $\lambda^l = (d^T \psi/C) \lambda^l \Rightarrow \mathcal{P}_r(\lambda) = 0$ 

with the *reduced polynomial* of degree l = 352:

$$\mathcal{P}_r(\lambda) = \lambda^l - \sum_{j=0}^{l-1} \lambda^{l-1-j} c_j = 0 \quad , \quad c_j = d^T S_0^j e/N \; .$$

 $\Rightarrow$  at most l = 352 eigenvalues  $\lambda \neq 0$  which can be numerically determined as the zeros of  $\mathcal{P}_r(\lambda)$ . However: still numerical problems:

• 
$$c_{l-1} \approx 3.6 \times 10^{-352}$$

- alternate sign problem with a strong loss of significance.
- big sensitivity of eigenvalues on  $c_j$

## Solution:

Using the multi precision library GMP with 256 binary digits the zeros of  $\mathcal{P}_r(\lambda)$  can be determined with accuracy  $\sim 10^{-18}$ .

Furthermore the Arnoldi method can also be implemented with higher and precision.

 $\underline{\rm red\ crosses}:$  zeros of  $\mathcal{P}_r(\lambda)$  from 256 binary  $_{\rm \tiny -0.5}$  digits calculation

<u>blue squares</u>: eigenvalues from Arnoldi method  $_{0.5}$ with 52, 256, 512, 1280 binary digits. In the last  $_{0}$  case:  $\Rightarrow$  break off at  $n_A = 352$  with vanishing \_ $_{0.5}$  coupling element.



## **Full Physical Review network**

Accurate eigenvalue spectrum for the full Physical Review network by a new rational interpolation method (left) and the HP Arnoldi method

(right):



### **Fractal Weyl law**



 $N_{\lambda}$  = number of complex eigenvalues with  $\lambda_c \leq |\lambda| \leq 1$ .  $N_t$  = reduced network size of Physical Review at time t.

$$N_{\lambda} = a N_t^b$$

# Perron-Frobenius matrix for chaotic maps

A new variant of the *Ulam Method* to construct the *Perron-Frobenius matrix* for the case of a mixed phase space:

Subdivide phase space in square cells of size  $M^{-1}$  and iterate a classical trajectory  $(t \sim 10^{11} - 10^{12})$  and attribute a new number to each new cell which is entered. At the same time count the number of transitions from cell i to cell  $j \Rightarrow N \times N$ -PF-Matrix (N=number of non-empty cells) by:

$$G_{ji} = \frac{n_{ji}}{\sum_{l} n_{li}}$$

Example: Chirikov map at  $k=k_c=0.971635406$  with M=10.



## **Eigenvalues**

for  $M=10,\,t=10^6$  and N=35





Phase space representation of the eigenvector for  $\lambda_0 = 1$ .









 $\lambda_0 = 1, M = 1600, N = 494964, n_A = 3000$ 







## **Extrapolation of eigenvalues**



## Absorption for p < 0.05

Chirikov map



Red, green (left): Survial Monte-Carlo Method Blue (left): Data of Weiss et al. PRL **89**, 239401 (2002) and Chirikov et al. PRL **89**, 239402 (2002).

# References

- 1. D. L. Shepelyansky *Fractal Weyl law for quantum fractal eigenstates*, Phys. Rev. E **77**, p.015202(R) (2008).
- L. Ermann and D. L. Shepelyansky, *Ulam method and fractal Weyl law for Perron-Frobenius operators*, Eur. Phys. J. B 75, 299 (2010).
- 3. K. M. Frahm and D. L. Shepelyansky, *Ulam method for the Chirikov standard map*, Eur. Phys. J. B **76**, 57 (2010).
- K. M. Frahm, B. Georgeot and D. L. Shepelyansky, *Universal emergence of PageRank*, J. Phys. A: Math. Theor. 44, 465101 (2011).
- 5. K. M. Frahm, A. D. Chepelianskii and D. L. Shepelyansky, *PageRank of integers*, arxiv:1205.6343[cs.IR] (2012).

- K. M. Frahm and D. L. Shepelyansky, *Google matrix of Twitter*, Eur. Phys. J. B 85, 355 (2012).
- L. Ermann, K. M. Frahm and D. L. Shepelyansky, Spectral properties of Google matrix of Wikipedia and other networks, Eur. Phys. J. B 86, 193 (2013).
- K. M. Frahm and D. L. Shepelyansky, Poincaré recurrences and Ulam method for the Chirikov standard map, Eur. Phys. J. B 86, 322 (2013).
- K. M. Frahm, and D. L. Shepelyansky, Poisson statistics of PageRank probabilities of Twitter and Wikipedia networks, Eur. Phys. J. B, 87, 93 (2014).
- K. M. Frahm, Y. H. Eom, and D. L. Shepelyansky, Google matrix of the citation network of Physical Review, Phys. Rev. E 89, 052814 (2014).
- L. Ermann, K. M. Frahm and D. L. Shepelyansky, Google matrix analysis of directed networks, submitted to Rev. Mod. Phys. July 25, (2014) (arXiv:1409.0428 [physics.soc-ph] preprint)