



RECHERCHE & TRAITEMENT DU LANGAGE

D'INFORMATION AUTOMATIQUE

RI - Recherche web

2022-23

Benjamin Piwowarski / Laure Soulier



Machine Learning &
Deep Learning for
Information Access

INTRODUCTION

Analyse de liens

- Popularisée par Google avec PageRank
- Actuellement une composante parmi beaucoup d'autres des moteurs de recherche
- De l'ordre de 400 caractéristiques prises en compte
- Cours : 2 algorithmes historiques
 - PageRank (Brin & Page 1998)
 - HITS (Kleinberg 1998)
 - Très nombreuses variantes

E.g. trustrank

- Le web est vu comme un graphe orienté
- Les liens sont porteurs d'information
 - Un lien entre pages indique une relation de pertinence
Un lien est un indicateur de qualité
 - Le texte d'un lien résume la page cible
L'indexation d'une page doit prendre en compte les liens vers cette page (contexte)

PAGERANK

PageRank in 1 slide

- Principe général
 - Popularized by google
 - Assign an authority score for each web page
 - Using only the structure of the web graph (query independent)
 - Now one of the many components used for computing page scores in Google S.E.
- Intuition
 - Assign higher scores to pages with many in-links from authoritative pages with few out-links
- Modèle
 - Random surfer model : Stationary distribution of a Markov Chain
 - Principal eigenvector of a linear system

Notations

- Graphe orienté $G = (V, E)$
- A matrice d'adjacence

$$a_{ij} = \begin{cases} 1 & \text{s'il y a un lien de } i \text{ vers } j \\ 0 & \text{sinon} \end{cases}$$

- Le nombre de liens entrants pour une page i est

$$d_i = \sum_j a_{ij}$$

- P matrice de transition

$$P = \left(\frac{a_{ij}}{d_i} \right)_{ij}$$

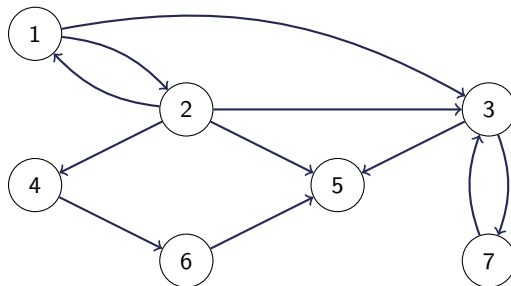
avec P_{ij} la probabilité de transition, i.e. d'aller de j à i . On a :

$$\sum_j p_{ij} = 1$$

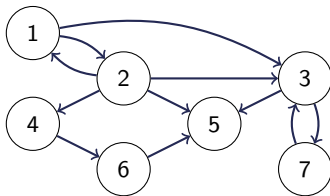
i.e.

$$p_{ij} \stackrel{\text{def}}{=} p(\text{aller sur } j | \text{point de départ } i)$$

Graphe

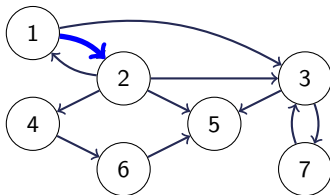


Représentation matricielle



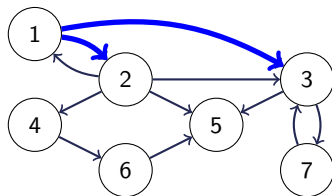
$$A = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Représentation matricielle



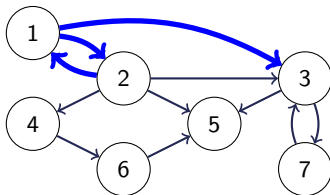
$$A = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Représentation matricielle



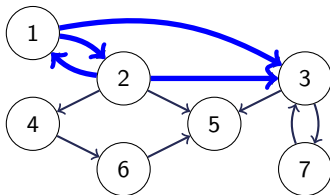
$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Représentation matricielle



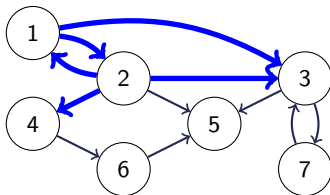
$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Représentation matricielle



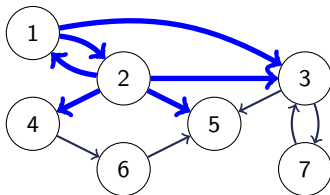
$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Représentation matricielle



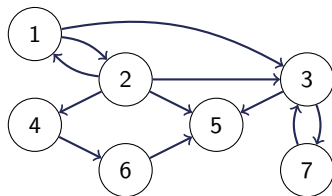
$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Représentation matricielle



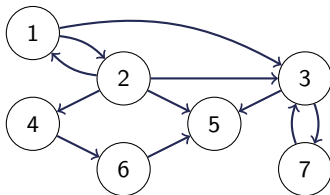
$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Représentation matricielle



$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

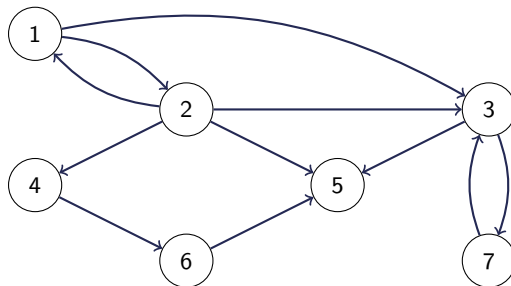
Représentation matricielle



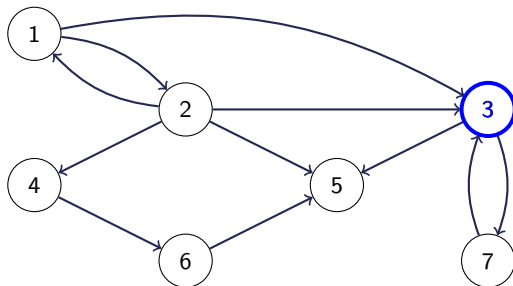
$$A = \begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$P = \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 \\ \frac{1}{4} & 0 & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{2} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{2} \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

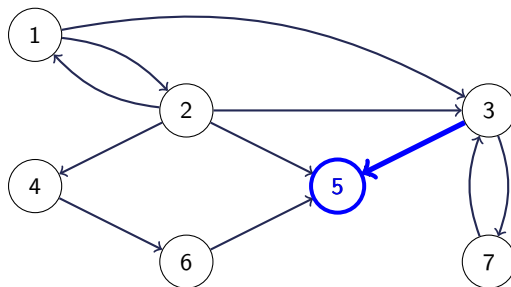
Surfer stochastique



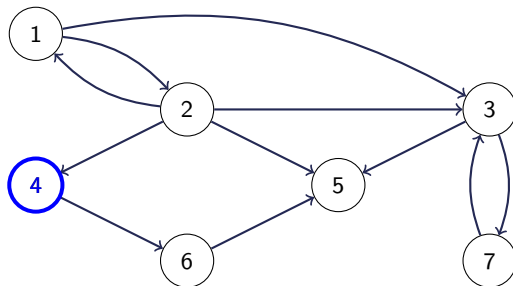
Surfer stochastique



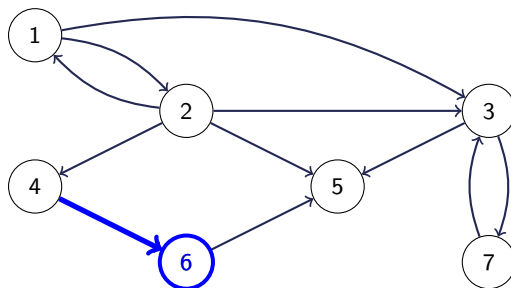
Surfer stochastique



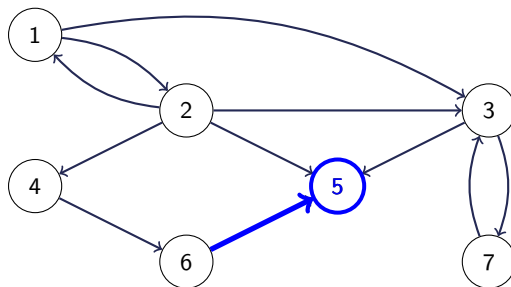
Surfer stochastique



Surfer stochastique



Surfer stochastique



Distribution stationnaire s_j

$$s_j = \underbrace{\sum_i p_{ij} s_i}_{\text{propagation de l'importance}} \quad (1)$$

- s_j correspond à la probabilité que la page i soit importante $\sum_i s_i = 1$

Distribution stationnaire \mathbf{s}_j

$$\mathbf{s}_j = \underbrace{d \sum_i p_{ij} \mathbf{s}_i}_{\text{propagation de l'importance}} + \underbrace{(1 - d) a_j}_{\text{a priori}} \quad (1)$$

- d est le facteur d'amortissement (*damping factor*). Les valeurs typiques de d sont autour de 0.8
- a_i correspond à la probabilité que la page i soit importante *a priori*
 $\sum_i a_i = 1$

Vision matricielle

- Version initiale

$$\mathbf{s}_j = \underbrace{d \sum_i \mathbf{p}_{ij} \mathbf{s}_i}_{\text{propagation de l'importance}} + \underbrace{(1 - d) \mathbf{a}_j}_{\text{a priori}}$$

Vision matricielle

- Version initiale

$$s_j = \underbrace{d \sum_i p_{ij} s_i}_{\text{propagation de l'importance}} + \underbrace{(1-d) a_j}_{\text{a priori}}$$

- Version matricielle

$$\mathbf{s} = d\mathbf{sP} + (1-d)\mathbf{a}$$

Un peu d'algèbre linéaire...

$$\mathbf{s} = d\mathbf{sP} + (1 - d)\mathbf{a} = \mathbf{s}(d\mathbf{P} + (1 - d)\mathbf{E}) = \mathbf{sP}'$$

Un peu d'algèbre linéaire...

$$\mathbf{s} = d\mathbf{sP} + (1 - d)\mathbf{a} = \mathbf{s}(d\mathbf{P} + (1 - d)\mathbf{E}) = \mathbf{sP}'$$

- Que vaut \mathbf{s} ?

Un peu d'algèbre linéaire...

$$\mathbf{s} = d\mathbf{sP} + (1 - d)\mathbf{a} = \mathbf{s}(d\mathbf{P} + (1 - d)\mathbf{E}) = \mathbf{sP}'$$

- Que vaut \mathbf{s} ?
- Rappel sur les valeurs propres : $A\mathbf{X} = \lambda\mathbf{X}$ (λ : valeur propre et \mathbf{X} : vecteur propre)

Un peu d'algèbre linéaire...

$$\mathbf{s} = d\mathbf{sP} + (1 - d)\mathbf{a} = \mathbf{s}(d\mathbf{P} + (1 - d)\mathbf{E}) = \mathbf{sP}'$$

- Que vaut \mathbf{s} ?
- Rappel sur les valeurs propres : $\mathbf{AX} = \lambda\mathbf{X}$ (λ : valeur propre et \mathbf{X} : vecteur propre)
- \mathbf{s} est le vecteur propre associé à la valeur propre 1

Un peu d'algèbre linéaire...

$$\mathbf{s} = d\mathbf{sP} + (1 - d)\mathbf{a} = \mathbf{s}(d\mathbf{P} + (1 - d)\mathbf{E}) = \mathbf{sP}'$$

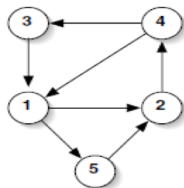
- Que vaut \mathbf{s} ?
- Rappel sur les valeurs propres : $\mathbf{AX} = \lambda\mathbf{X}$ (λ : valeur propre et \mathbf{X} : vecteur propre)
- \mathbf{s} est le vecteur propre associé à la valeur propre 1
- Il existe une solution unique

Explication théorique

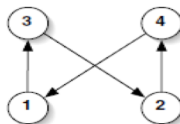
A est une matrice irréductible : correspond à un graphe fortement connecté

- A square matrix $A_{n \times n}$ is **non negative** if $a_{ij} \geq 0$
 - Notation $A \geq 0$
 - Example: graph incidence matrix
- $A_{n \times n}$ is **positive** if $a_{ij} > 0$
 - Notation $A > 0$
- $A_{n \times n}$ is **irreducible** if
 - $\forall i, j, \exists t \in \mathbb{N} / (A^t)_{ij} > 0$
 - If A is a graph incidence matrix, this means that G is strongly connected
 - There is a path between any pair of vertices
- $A_{n \times n}$ is **primitive** if $\exists t \in \mathbb{N} / A^t > 0$
 - A primitive matrix is irreducible
 - Converse is false

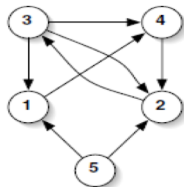
Explication théorique



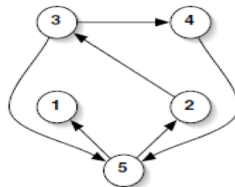
(a)



(b)



(c)

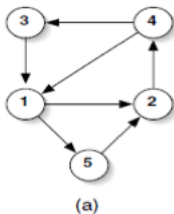


(d)

Figure 5.1 Graphs with different types of incidence matrices. (a) is primitive, (b) is irreducible (with period 4) but not primitive, (c) and (d) are reducible.

Explication théorique

Exemples (Baldi et al. 2003)



$$\begin{pmatrix} 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{pmatrix}^7 = \begin{pmatrix} 4 & 2 & 3 & 3 & 1 \\ 3 & 4 & 1 & 1 & 3 \\ 1 & 3 & 1 & 3 & 1 \\ 2 & 6 & 1 & 4 & 3 \\ 3 & 1 & 2 & 1 & 1 \end{pmatrix}$$

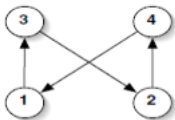
Matrice Primitive

« Il existe une puissance pour laquelle tous les éléments sont strictement positifs »

En termes de graphe : On peut naviguer entre tous les noeuds

Explication théorique

Exemples (Baldi et al. 2003)



Irréductible :

« Il existe un exposant qui permet pour tout élément d'avoir une valeur de 1 »

En termes de graphe : On peut naviguer entre tous les noeuds

$$\begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}^1 = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

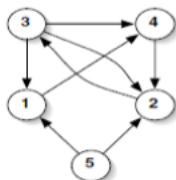
$$\begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}^2 = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

≡ ≡

$$\begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}^3 = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

$$\begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}^4 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Exemples (Baldi et al. 2003)



(c)

$$\begin{pmatrix} 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \end{pmatrix}^{10} = \begin{pmatrix} 4 & 8 & 5 & 6 & 0 \\ 8 & 17 & 11 & 13 & 0 \\ 11 & 24 & 17 & 19 & 0 \\ 5 & 11 & 8 & 9 & 0 \\ 7 & 16 & 12 & 13 & 0 \end{pmatrix}$$

Réductible :

« On a toujours un bloc de 0, quelque que soit la puissance »

En termes de graphe : On ne peut naviguer entre tous les noeuds

Définitions (Matrice irréductible)

- A est irréductible si $\exists n \mathbf{A}^n > 0$
- Correspond à un graphe fortement connecté
- Si $\mathbf{A}_{n \times n}$ est une matrice non-négative irréductible et apériodique
 - A a une valeur propre réelle λ_1 et positive supérieure unique telle que $\forall j, |\lambda_1| > |\lambda_j|$
 - Le vecteur propre s associé à λ_1 est strictement positif
- Dans notre cas, la valeur propre maximum est ?

Définitions (Matrice irréductible)

- A est irréductible si $\exists n \mathbf{A}^n > 0$
- Correspond à un graphe fortement connecté
- Si $\mathbf{A}_{n \times n}$ est une matrice non-négative irréductible et apériodique
 - A a une valeur propre réelle λ_1 et positive supérieure unique telle que $\forall j, |\lambda_1| > |\lambda_j|$
 - Le vecteur propre s associé à λ_1 est strictement positif
- Dans notre cas, la valeur propre maximum est 1

Preuve de convergence

- Méthodes des puissances

$$A = U\Sigma U^t \rightarrow A^n = U\Sigma^n U^t$$

$$\Sigma^n \xrightarrow{\infty} \begin{pmatrix} 1 & & & \\ & 0 & & \\ & & 0 & \\ & & & 0 \end{pmatrix}$$

Revient à sélectionner le premier vecteur propre

En pratique...

- Version itérative :
 - Initialiser s aléatoirement
 - Répéter pour chaque noeud

$$s_j = d \sum_i p_{ij} s_i + (1 - d) a_j$$

- Dans les deux cas, jusqu'à ce que $|s^{(t+1)} - s^{(t)}| < \epsilon$
- Rapidité de la convergence : géométrique avec un ratio λ_1/λ_2

Démonstration

$$\mathbf{A} = \begin{pmatrix}
 . & . & 1.0 & 1.0 & 1.0 & . & 1.0 & . & 1.0 & . \\
 . & . & . & 1.0 & 1.0 & 1.0 & . & . & . & . \\
 1.0 & . & . & 1.0 & 1.0 & 1.0 & . & . & . & . \\
 1.0 & . & 1.0 & . & 1.0 & 1.0 & 1.0 & 1.0 & 1.0 & 1.0 \\
 . & . & 1.0 & 1.0 & . & 1.0 & . & 1.0 & 1.0 & . \\
 1.0 & 1.0 & 1.0 & . & 1.0 & . & 1.0 & . & 1.0 & 1.0 \\
 . & 1.0 & . & . & 1.0 & 1.0 & . & 1.0 & . & . \\
 1.0 & 1.0 & . & 1.0 & 1.0 & 1.0 & . & . & 1.0 & . \\
 . & 1.0 & 1.0 & 1.0 & 1.0 & 1.0 & . & . & . & . \\
 1.0 & 1.0 & 1.0 & 1.0 & . & 1.0 & 1.0 & . & . & .
 \end{pmatrix}$$

Démonstration

$$\mathbf{P} = \begin{pmatrix}
 . & . & 0.2 & 0.2 & 0.2 & . & 0.2 & . & 0.2 & . \\
 . & . & . & 0.33 & 0.33 & 0.33 & . & . & . & . \\
 0.25 & . & . & 0.25 & 0.25 & 0.25 & . & . & . & . \\
 0.12 & . & 0.12 & . & 0.12 & 0.12 & 0.12 & 0.12 & 0.12 & 0.12 \\
 . & . & 0.2 & 0.2 & . & 0.2 & . & 0.2 & 0.2 & . \\
 0.14 & 0.14 & 0.14 & . & 0.14 & . & 0.14 & . & 0.14 & 0.14 \\
 . & 0.25 & . & . & 0.25 & 0.25 & . & 0.25 & . & . \\
 0.17 & 0.17 & . & 0.17 & 0.17 & 0.17 & . & . & 0.17 & . \\
 . & 0.2 & 0.2 & 0.2 & 0.2 & 0.2 & . & . & . & . \\
 0.17 & 0.17 & 0.17 & 0.17 & . & 0.17 & 0.17 & . & . & .
 \end{pmatrix}$$

Démonstration

Iteration 1, 2, 5, 10

$$\mathbf{x}_1^\top = \begin{pmatrix} 0.07 \\ 0.078 \\ 0.13 \\ 0.17 \\ 0.14 \\ 0.15 \\ 0.07 \\ 0.062 \\ 0.12 \\ 0.015 \end{pmatrix}, \mathbf{x}_2^\top = \begin{pmatrix} 0.088 \\ 0.075 \\ 0.11 \\ 0.14 \\ 0.17 \\ 0.16 \\ 0.06 \\ 0.066 \\ 0.095 \\ 0.043 \end{pmatrix}, \mathbf{x}_5^\top = \begin{pmatrix} 0.086 \\ 0.075 \\ 0.12 \\ 0.14 \\ 0.16 \\ 0.16 \\ 0.064 \\ 0.065 \\ 0.1 \\ 0.04 \end{pmatrix}, \mathbf{x}_{10}^\top = \begin{pmatrix} 0.086 \\ 0.076 \\ 0.12 \\ 0.14 \\ 0.16 \\ 0.16 \\ 0.064 \\ 0.065 \\ 0.1 \\ 0.04 \end{pmatrix}$$

HITS

Hubs : page de pointeurs



Cellarer search

Top 200 wine sites

Here is a directory of websites and blogs about wine. The table **only shows the top 200** sites out of **a list of 500**. Other pages explain the ranking criteria: [website valuation](#), [Google PageRank](#), [monthly traffic](#). Please [comment](#) on the results.

As a complement the Cellarer search engine shows you lesser-known but interesting authors. Take the tour: [Japanese recipes](#), [local](#), [Bordeaux](#).

Website	Cellarer rank	Valuation	PageRank	Monthly traffic
Smooth	1	\$244 803	6	332962
Wine Spectator	2	\$88 238	7	87712
Wine Enthusiast	3	\$29 072	6	30936
Dr. Wine	4	\$28 690	6	30403
Cork'd	5	\$27 280	6	28429
Wine lovers page	6	\$25 812	6	26374
The Winedoctor	7	\$25 384	6	25775
eRobert Parker	8	\$24 817	6	24981
Decanter	9	\$22 895	6	22291

- Points good authority pages

Authorities : pages de références thématiques



- Pointed by good hub pages
- Référence importante pour un thème

Hubs et Authorities

- Pour une page i
 - Score hub = somme des score des *authorities* des pages pointées par i
 - Score authority = somme des scores des hubs des pages qui pointent vers i
- Formellement,

$$h_i = \sum_{i \rightarrow j} a_j$$

$$a_i = \sum_{j \rightarrow i} h_j$$

Algorithme

- Itératif
 - On calcule a et h à partir des valeurs estimées à l'itération précédente
 - il faut normaliser a et h à chaque étape pour que cela fonctionne (norme L2 égale à 1)
- Version algèbre linéaire

$$\lambda_h \mathbf{h} = \mathbf{h} \mathbf{P} \mathbf{P}^\top$$

$$\lambda_a \mathbf{a} = \mathbf{a} \mathbf{P}^\top \mathbf{P}$$