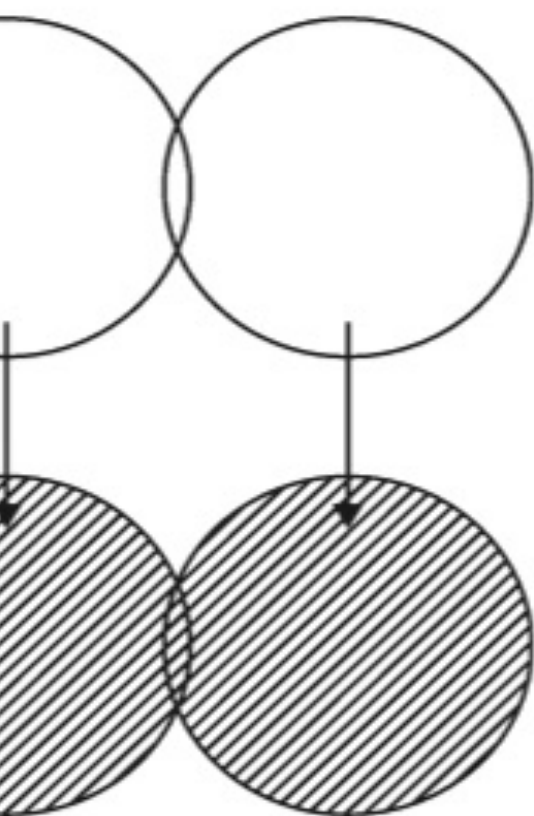


M1 – Algèbre linéaire appliquée

Google PageRank

D e c e m b e r , 2 0 2 3



Lidia Nievas
Malang Badiane

Projet d'algèbre lineaire

MALANG BADIANE

LIDIA NIEVAS DUENAS

December 26, 2023

1 Introduction

Le moteur de recherche Google a été mis au point à la fin des années 90 par S.Brin et L.Page. À l'issue d'une requête, il affiche les réponses dans un ordre reposant sur un indice de popularité des pages web : le PageRank. Ce procédé a permis à Google de devancer la concurrence. Le but du projet sera d'implémenter le calcul de l'indice : il s'agit d'un calcul de vecteur propre par la méthode de puissance itérée. On supposera que le web est constitué des quatre pages données ci-dessous, les caractères en gras indiquent les hyperliens :

1.1 P1 : Matrice des hyperliens

La matrice des hyperliens code les liens entre les pages du web. la **matrice de Google** est un employé par le célèbre moteur de recherche repose sur cette matrice.

1.2 P2 : Matrice stochastique

Une matrice stochastique est une matrice carrée dont les coefficients sont positifs et dont la somme des coefficients vaut 1 sur chaque ligne. La **matrice de Google** est un exemple célèbre de matrice stochastique, dont le **vecteur propre** de la valeur propre dominante fournit le PageRank.

1.3 P3 : Matrice de Google

La matrice de Google est une **matrice stochastique** déduite de la matrice des hyperliens de l'ensemble des pages Web. Le vecteur PageRank de classement des pages par popularité est un **vecteur propre** associé à la plus grande valeur propre de cette matrice.

1.4 P4 : Vecteur propre

Si \mathbf{A} est une matrice réelle $n \times n$, on dit que le vecteur $v \in \mathbb{R}^n$ est un vecteur propre de \mathbf{A} si $v \neq 0$ et s'il existe $\lambda \in \mathbb{C}$ tel que $\mathbf{A}v = \lambda v$.

Question 1

Sur cet exemple, donnons les matrices d'hyperlien H , stochastique S , puis celle de G en prenant $\alpha = 0.85$, et le vecteur v uniforme :

$$v = \frac{1e}{n}$$

On définit la matrice d'hyperlien H par : $H_{i,j} = \begin{cases} \frac{1}{p_{i,j}} & \text{si } P_i \text{ possède un lien qui pointe vers la page } P_j \\ 0 & \text{sinon} \end{cases}$

Ainsi on aura dans notre cas :

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

On déduit la matrice stochastique S :

$$S = H + \frac{1}{n} \cdot e^T$$

$$S = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 1/2 & 1/2 \\ 0 & 1/2 & 0 & 1/2 \\ 1/4 & 1/4 & 1/4 & 1/4 \end{pmatrix}$$

On définit la matrice de google par :

$$G = \alpha S + (1 - \alpha) e v^T$$

Ce qui nous donne pour $\alpha = 0.85$ et $v = \frac{1}{4}e$

$$G = \begin{pmatrix} 0.0375 & 0.0375 & 0.0375 & 0.0375 \\ 0.0375 & 0.0375 & 0.4625 & 0.4625 \\ 0.0375 & 0.4625 & 0.0375 & 0.4625 \\ 1/4 & 1/4 & 1/4 & 1/4 \end{pmatrix}$$

Code sur Matlab:

```

1 %Exercice 1
2 disp('Exercice 1')
3 n=4;
4 %Matrice H
5 H=[0,0,1,0;0,0,1/2,1/2;0,1/2,0,1/2;0,0,0,0];
6 disp('H=')
7 disp(H)
8 %Matrice S
9 S=[0,0,1,0;0,0,1/2,1/2;0,1/2,0,1/2;1/4,1/4,1/4,1/4];
10 disp('S=')
11 disp(S)
12 %Matrice G
13 alpha=0.85;
14 e=[1;1;1;1];
15 v=(1/n)*e;
16 G=alpha*S+(1-alpha)*e*v';
17 disp('G=')
18 disp(G)
19
20

```

Sortie d'écran:

```

Exercice 1
H=
    0         0    1.0000         0
    0         0    0.5000    0.5000
    0    0.5000         0    0.5000
    0         0         0         0

S=
    0         0    1.0000         0
    0         0    0.5000    0.5000
    0    0.5000         0    0.5000
    0.2500    0.2500    0.2500    0.2500

G=
    0.0375    0.0375    0.8875    0.0375
    0.0375    0.0375    0.4625    0.4625
    0.0375    0.4625    0.0375    0.4625
    0.2500    0.2500    0.2500    0.2500
fx

```

2 L'algorithme de la puissance itérée :

2.1 Question 2

- Le pseudo-code :

Entrés : matrice stochastique $A < 0$ de taille $n \times n$, vecteur ligne stochastique $x_0 \geq 0$ d'initialisation, tolérance $\epsilon > 0$ le résidu, borne $nmax$ sur le nombre d'itération.

Poser $z = x_0 A$

Calculer le résidu : $\eta = \|z - x\|_1$

mettre à jour : $x = z, k = 1$

Tant Que : $\eta > \epsilon$ et $k < nmax$

Poser $z = xA$,

calculer le 'résidu' : $\eta = \|z - x\|_1$

Mettre à jour : $x = z, k = k + 1$

Fin TantQue

Sorties : vecteur ligne stochastique $x > 0$ tel que $xA = x$, et nombre k d'itérations.

Le code sur Matlab :

```
223 %Fonction exercice 2
224 function [x,k]=puissanceiteree(A,x0,epsilon,nmax)
225     %Entrée: matrice stochastique A
226     %vecteur ligne stochastique x0 d'initialisation
227     %entier epsilon, la tolérance
228     %entier nmax, nombre max d'itérations
229     %Sortie: le PageRank pi
230     %entier k, le nombre d'itérations
231     z=x0*A;
232     eta=norm(z-x0,1);
233     x=z;
234     k=1;
235     while (eta>epsilon && k<nmax)
236         z=x*A;
237         eta=norm(z-x,1);
238         x=z;
239         k=k+1;
240     end
241 end
242
243
```

2.2 Question 3

On calcule le pagerank des pages :

```
20
21 %Exercice 3
22 disp('Exercice 3')
23 %Définition des variables
24 pi0=[0.25,0.25,0.25,0.25];
25 epsilon=0.01;
26 nmax=100;
27 %Calcul du PageRank avec puissanceiteree
28 [pi,k]=puissanceiteree(G,pi0,epsilon,nmax);
29 disp('pi=')
30 disp(pi)
31 disp('k=')
32 disp(k)
33
34
```

On obtient le résultat suivant :

```
Exercice 3
pi=
    0.1104    0.2413    0.3054    0.3428

k=
    6
```

On constate que la méthode converge en six itérations et le vecteur du pagerank est : $\pi = (0.1104, 0.2413, 0.3054, 0.3428)$
On constate que la page P4 est la plus populaire suivie de P3, P2 et P1

2.3 Question 4

On teste maintenant différents vecteurs de personnalisation :

```
35 %Exercice 4
36 disp('Exercice 4')
37 %Définition des variables
38 alpha=0.85;
39 e=[1;1;1;1];
40 v1=[0.1,0.4,0.1,0.4];
41 v2=[0.02,0.48,0.02,0.48];
42 G1=alpha*S+(1-alpha)*e*v1;
43 G2=alpha*S+(1-alpha)*e*v2;
44 pi0=[0.25,0.25,0.25,0.25];
45 epsilon=0.01;
46 nmax=100;
47 %Calcul du PageRank avec puissanceiterée
48 [pi1,k1]=puissanceiterée(G1,pi0,epsilon,nmax);
49 disp('pi1=')
50 disp(pi1)
51 disp('k1=')
52 disp(k1)
53 [pi2,k2]=puissanceiterée(G2,pi0,epsilon,nmax);
54 disp('pi2=')
55 disp(pi2)
56 disp('k2=')
57 disp(k2)
58
```

Sortie d'écran:

```
Exercice 4
pi1=
    0.0932    0.2586    0.2808    0.3674

k1=
    6

pi2=
    0.0839    0.2678    0.2677    0.3806

k2=
    6
```

On constate que l'indice de popularité de la page 4 est le plus grand dans les trois cas. Donc la page la plus populaire est la page 4. le vecteur de personnalisation v est choisi pour gonfler l'indice de popularité des pages. Lorsqu'on gonfle une composante du vecteur v on augmente son indice de popularité : par exemple en modifiant v_4 au lieu de 0.4 on met 0.48, l'indice de popularité de la page 4 est passée de 0.3674 à 0.3806.

3 Le PageRank

3.1 Question 5

On donne le code du pageRank :

```
59
60 %Exercice 5
61 disp('Exercice 5')
62 %Définition des variables
63 v1=[0.1;0.4;0.1;0.4];
64 v2=[0.02;0.48;0.02;0.48];
65 %Calcul du PageRank avec PageRank
66 [pi1,k1]=PageRank(H,v1,alpha,pi0,epsilon,nmax);
67 disp('pi1=')
68 disp(pi1)
69 disp('k1=')
70 disp(k1)
71 [pi2,k2]=PageRank(H,v2,alpha,pi0,epsilon,nmax);
72 disp('pi2=')
73 disp(pi2)
74 disp('k2=')
75 disp(k2)
76
193 %Fonction exercice 5
194 function[pi,k]=PageRank(H,v,alpha,pi0,epsilon,nmax)
195 %Entrée: matrice des hyperliens H
196 %Vecteur stochastique de personnalisation v
197 %entier alpha, le parametre
198 %vecteur ligne stochastique pi0 d'initialisation
199 %entier epsilon, la tolérance
200 %entier nmax, nombre max d'itérations
201 %Sortie: le PageRank pi
202 %entier k, le nombre d'itérations
203 n=size(v,1);
204 e=ones(1,n)';
205 a=zeros(1,n)';
206 for i=1:n
207     if(H(i,:)~=zeros(1,n))
208         a(i)=1;
209     end
210 end
211 z=alpha*(pi0*sparse(H)+(pi0*a)*e'/n)+(1-alpha)*v';
212 eta=norm(z-pi0,1);
213 pi=z;
214 k=1;
215 while (eta>epsilon && k<nmax)
216     z=alpha*(pi*sparse(H)+(pi*a)*e'/n)+(1-alpha)*v';
217     eta=norm(z-pi,1);
218     pi=z;
219     k=k+1;
220 end
221 end
222
```

Sortie d'écran:

```
Exercice 5
pi1=
    0.0932    0.2586    0.2808    0.3674

k1=
     6

pi2=
    0.0839    0.2678    0.2677    0.3806

k2=
     6
```

3.2 Question 6

La fonction genererH :

```
76
77 %Exercice 6
78 disp('Exercice 6')
79 [H]=genererH(5,3);
80 disp('H=')
81 disp(full(H))
82
83
243
244 %Fonction exercice 6
245 function[H]=genererH(n,m)
246     %Entrée: entier n, nombre des pages
247     %entier m, nombre des hyperliens
248     %Sortie: matrice des hyperliens H
249     if n<=m
250         disp("Erreur")
251     end
252     H=sparse(n,n);
253     for i=1:n
254         nbre_hyperliens=(m+1)*rand();
255         nombre_hyperliens=floor(nbre_hyperliens);
256         if nombre_hyperliens~=0
257             count=1;
258             while count<=nombre_hyperliens
259                 index=1+n*rand();
260                 indice=floor(index);
261                 if (H(i,indice)==sparse(1,1))
262                     H(i,indice)=1/nombre_hyperliens;
263                     count=count+1;
264                 end
265             end
266         end
267     end
268 end
269
```

Sortie d'èren:

```
Exercice 6
H=
    0         0    0.5000         0    0.5000
0.3333    0.3333         0    0.3333         0
1.0000         0         0         0         0
    0    1.0000         0         0         0
    0         0         0         0         0
```

3.3 Question 7

Comparaison du temps des calculs du pageRank et la puissance itérée:

Pour cette question nous executons les deux fonctions(pageRank,puissanceiterree) et à chaque fois nous conservons leurs temps d'exécution des deux fonctions.

```

untitled IP_alg.m +
84 %Exercice 7
85 disp('Exercice 7')
86 n=10000;
87 [H]=genererH(n,50);
88
89 tic
90 %Définition des variables pour écrire G
91 a=zeros(1,n);
92 %Remplissez le vecteur 'a' avec des 1 là où la matrice 'H' a une ligne pleine de 0
93 for i=1:n
94     if(H(i,:)==zeros(1,n))
95         a(i)=1;
96     end
97 end
98 nmax=100;
99 alpha=0.5;
100 pi0=ones(1,n)/n;
101 S=H+(a/n)*ones(1,n);
102 G=alpha*S+(1-alpha)*ones(1,n)*pi0';
103 %Calcul du PageRank avec puissanceiterée
104 [pi,k]=puissanceiterée(G,pi0,0.001,nmax);
105 toc
106
107 tic
108 %Définition des variables
109 nmax=100;
110 alpha=0.5;
111 pi0=ones(1,n)/n;
112 w=ones(1,n)/n;
113 v=w';
114 %Calcul du PageRank avec PageRank
115 [pi,k]=PageRank(H,v,alpha,pi0,0.001,nmax);
116 toc

```

Sortie d'ère:

```

Exercice 7
Elapsed time is 12.182010 seconds.
Elapsed time is 6.966459 seconds.

```

Pour une matrice de taille 1000×1000 , On constate que la méthode du PageRank est moins couteuse au niveau de temps par rapport à la méthode de la puissance itérée.

3.4 Question 8

Nos ordinateurs n'arrivent pas executer les codes, par soucis d'insufisance de mémoire pour stocké les variables de type 100000×100000 .

Mais le code est:

```

118 %Exercice 8
119 disp('Exercice 8')
120 n=100000;
121 [H]=genererH(n,50);
122
123 tic
124 %Définition des variables pour calculer G
125 a=zeros(1,n);
126 %Remplissez le vecteur 'a' avec des 1 là où la matrice 'H' a une ligne pleine de 0
127 for i=1:n
128     if(H(i,:)==zeros(1,n))
129         a(i)=1;
130     end
131 end
132 nmax=100;
133 a=sparse(a);
134 alpha=0.5;
135 pi0=ones(1,n)/n;
136 pi0=sparse(pi0);
137 S=H+(a/n)*ones(1,n);
138 S=sparse(S);
139 G=alpha*S+(1-alpha)*ones(1,n)*pi0';
140 G=sparse(G);
141 %Calcul du PageRank avec puissanceiterée
142 [pi,k]=puissanceiterée(G,pi0,0.001,nmax);
143 toc
144
145 tic
146 %Définition des variables
147 nmax=100;
148 alpha=0.5;
149 pi0=ones(1,n)/n;
150 pi0=sparse(pi0);
151 w=ones(1,n)/n;
152 v=w';
153 v=sparse(v);
154 %Calcul du PageRank avec PageRank
155 [pi,k]=PageRank(H,v,alpha,pi0,0.001,nmax);
156 toc

```


3.5 Question 9

Honorer le contrat : Pour honorer le contrat, il suffira de manipuler le vecteur v de personnalisation qui permettra d'élever le rang des pages concernées dans le vecteur colonne des indices de popularité. Pour le marché qui nous est proposé, on va augmenter l'indice de popularité des 5 pages P6, . . . , P10. Pour ce faire, on va gonfler les composante du vecteur v entre la 6ème et la 10ème et diminuer les autres, tout en veillant à ce que v soit toujours unitaire.

```
158 %Exercice 9
159 disp('Exercice 9')
160 %Définition des variables
161 n=100000;
162 [H]=genererH(n,50);
163 nmax=100;
164 alpha=0.5;
165 pi0=ones(1,n)/n;
166 w=ones(1,n)/n;
167 v=w';
168 %Calcul du PageRank
169 [pi,k]=PageRank(H,v,alpha,pi0,0.001,nmax);
170 disp(pi(1,6:10))
171 %Calcul du PageRank avec les pages 6 à 10 avec plus de poids
172 v0=per_v(n)';
173 [pi,k]=PageRank(H,v0,alpha,pi0,0.001,nmax);
174 disp(pi(1,6:10))
175
176 %Fonction pour implementer les poids sur les pages 6 à 10
177 function[v0]=per_v(n)
178 %Entrée: un entier n, la taille du matrice
179 %Sortie: un vecteur v0
180 %Sortie: un vecteur v0
181 v0=(1/n)*ones(1,n);
182 for i=6:10
183     v0(i)=1.5*(1/n);
184 end
185 for i=1:5
186     v0(i)=(1-sum(v0(6:10)))/(n-5);
187 end
188 for i=11:n
189     v0(i)=(1-sum(v0(6:10)))/(n-5);
190 end
191 end
```

Appliqué à une matrice d'une taille 10000x10000, l'ordinateur n'a pas assez de mémoire pour exécuter le programme, mais un exemple de ce qui se passerait avec une matrice plus petite (de taille 15x15 dans ce cas) est le suivant :

Exercice 9

PageRank normal

0.0676	0.0892	0.0770	0.0512	0.0579
--------	--------	--------	--------	--------

PageRank avec les pages 6 à 10 avec plus de poids

0.0811	0.1065	0.0902	0.0660	0.0741
--------	--------	--------	--------	--------

4 Conclusion

Les moteurs de recherche sont d'énormes facteurs de pouvoir sur le Web, guidant les gens vers l'information et prestations de service.

Google est le moteur de recherche le plus performant de ces dernières années, ses résultats de recherche très complets et précis. Quand Google n'était qu'un projet de recherche précoce à Stanford, plusieurs articles ont été écrits décrivant les algorithmes sous-jacents.

L'algorithme dominant a été appelé PageRank et est toujours la clé pour fournir classements précis pour les résultats de recherche. Une fonctionnalité primordiale des moteurs de recherche de pages web consiste en un tri des résultats associés à une requête par ordre d'importance ou de pertinence.

Nous présentons un modèle permettant de définir une quantification de cette notion (Pagerank) a priori floue et des éléments de formalisation pour la résolution numérique du problème. On commence par une première approche naturelle non satisfaisante dans certains cas. Un raffinement de l'algorithme est donc introduit pour améliorer les résultats.