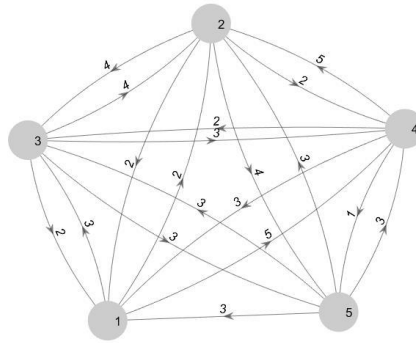


Projet SINF 1250 « Ranking de pages web : PageRank » (groupe 3)



Introduction : qu'est-ce que PageRank ?

PageRank est un algorithme de classement de page Internet qui fonctionne en comptant le nombre et la qualité des liens vers une page afin d'estimer l'importance du site en question. L'intuition est que les sites les plus importants sont plus susceptibles de recevoir plus de lien venant d'autres sites.

Matrice d'adjacence du graph orienté sur Moodle :

Pour rappel, dans le cas d'un graph orienté, au sein de la matrice d'adjacence, l'élément a_{ij} désigne le nombre d'arc d'origine i et d'extrémité j . Le terme à l'indice ij compte toujours le nombre de chemins allant de i vers j .

Dans le cadre de ce projet, en plus d'être orienté, le graph est pondéré ce qui nous amène à la matrice d'adjacence d'un graph pondéré. Il s'agit d'une matrice dans laquelle on note à la place du nombre de lien entre 2 nœud, la somme du poids des arcs entre 2 nœuds.

Matrice d'adjacence de notre graph :

$$\mathbf{A} = \begin{pmatrix} 0 & 2 & 3 & 5 & 0 \\ 1 & 0 & 4 & 2 & 4 \\ 2 & 4 & 0 & 3 & 3 \\ 3 & 5 & 2 & 0 & 1 \\ 3 & 3 & 3 & 3 & 0 \end{pmatrix}$$

Dans le cadre de notre programme python, nous avons donné en input la matrice d'adjacence sous forme d'un fichier muni de l'extension « .csv ». Chaque élément a_{ij} de la matrice est séparé par une virgule et les lignes sont représenté en faisant un retour à la ligne. Notre programme a ensuite lis ce fichier pour le convertir en matrice dans python. La lecture du fichier se fait au sein de la méthode `main()` du programme et utilise les méthodes de la librairie de python destinées à cet usage.

Comment calculer le vecteur de scores PageRank ?

Matrice de probabilité de transition :

Les conditions pour pouvoir déterminer un graph probabiliste est qu'il doit être orienté et pondéré.

La matrice de transition associé à un graph probabiliste d'ordre n est une matrice carrée d'ordre n . Pour trouver le facteur à la position a_{ij} de la matrice de transition, nous allons utiliser la formule du cours :

Formule du cours (slide 123) :

$$P(\text{page}(k+1) = i | \text{page}(k) = j) = \frac{w_{ji}}{w_j.}$$

$$w_{j.} = \sum_{i=1}^n w_{ji}$$

Cette formule décrit : la probabilité de passer d'une page j , en étant à un page i , est la probabilité de passer de la page j à i qui correspond à la pondération de l'arc de j vers i divisé par la somme des probabilités pour passer de la page j à toutes les autres pages i du graph qui correspond à la somme des pondérations de j vers tous les autres pages i . La probabilité de passer de la page j à toutes les autres par i du graph correspond au degré sortant du nœud j .

On réalise cette suite d'opération pour chaque nœud et on obtient la matrice de probabilité de transition où les lignes correspondent aux nœuds de départ et les colonnes aux nœuds d'arrivée.

Matrice de probabilité de transition :

$$\mathbf{P} = \begin{pmatrix} \frac{0}{10} & \frac{2}{10} & \frac{3}{10} & \frac{5}{10} & \frac{0}{10} \\ \frac{1}{11} & \frac{0}{11} & \frac{4}{11} & \frac{2}{11} & \frac{4}{11} \\ \frac{11}{12} & \frac{11}{12} & \frac{11}{12} & \frac{11}{12} & \frac{11}{12} \\ \frac{2}{12} & \frac{4}{12} & \frac{0}{12} & \frac{3}{12} & \frac{3}{12} \\ \frac{3}{12} & \frac{5}{12} & \frac{2}{12} & \frac{0}{12} & \frac{1}{12} \\ \frac{11}{12} & \frac{11}{12} & \frac{11}{12} & \frac{11}{12} & \frac{11}{12} \\ \frac{3}{12} & \frac{3}{12} & \frac{3}{12} & \frac{3}{12} & \frac{0}{12} \end{pmatrix}$$

L'avant dernière étape consiste à remplir la matrice Google en utilisant la formule :

$$G_{ij} = \alpha * P_{ij} + (1-\alpha) * 1/N$$

G_{ij} : l'élément à l'indice ij de la matrice Google

Alpha : le facteur de saut "dumping factor" qui nous est donné et qui vaut 0.9

P_{ij} : l'élément à l'indice ij de la matrice de transition

N : le nombre total de nœud de graph (qui correspond au nombre total de pages web)

On obtient :

$$G = \begin{pmatrix} \frac{1}{50} & \frac{9}{50} & \frac{27}{100} & \frac{9}{20} & 0 \\ \frac{9}{9} & \frac{1}{18} & \frac{18}{9} & \frac{9}{9} & \frac{18}{9} \\ \frac{110}{3} & \frac{50}{3} & \frac{55}{1} & \frac{55}{9} & \frac{55}{9} \\ \frac{20}{27} & \frac{10}{9} & \frac{50}{9} & \frac{40}{1} & \frac{40}{9} \\ \frac{110}{9} & \frac{22}{9} & \frac{55}{9} & \frac{50}{9} & \frac{110}{1} \\ \frac{40}{40} & \frac{40}{40} & \frac{40}{40} & \frac{40}{40} & \frac{50}{50} \end{pmatrix}$$

$$G^T = \begin{pmatrix} \frac{1}{50} & \frac{9}{110} & \frac{3}{20} & \frac{27}{110} & \frac{9}{40} \\ \frac{9}{9} & \frac{1}{1} & \frac{3}{3} & \frac{9}{9} & \frac{9}{9} \\ \frac{50}{27} & \frac{50}{18} & \frac{10}{1} & \frac{22}{9} & \frac{40}{9} \\ \frac{100}{9} & \frac{55}{9} & \frac{50}{9} & \frac{55}{1} & \frac{40}{9} \\ \frac{20}{9} & \frac{55}{18} & \frac{40}{9} & \frac{50}{9} & \frac{40}{1} \\ 0 & \frac{18}{55} & \frac{9}{40} & \frac{9}{110} & \frac{1}{50} \end{pmatrix}$$

Explications de la formule de la matrice de Google :

La matrice G décrit ceci : La plupart du temps, quelqu'un qui surf sur le net va suivre un lien sur une page : de la page i la personne va suivre les autres liens qui mène au voisin de i. Un petit pourcent du temps, la personne va quitter la page pour en rejoindre une différente, il va s'y « téléporter ». Le facteur P reflète la probabilité que la personne quitte la page actuelle et se « téléporte » vers une nouvelle. Comme il/elle peut se téléporter vers n'importe quelle page, chaque page à 1/N probabilité d'être choisie.

La dernière étape consiste à utiliser la power méthode afin de déterminer le vecteur contenant les page rank de chaque page internet.

$$v0 * G^T = v1$$

v0 : le vecteur initiale de l'itération

Transposée(G) : matrice de Google transposée

v1 : le nouveau vecteur qui servira pour l'itération suivante

Le vecteur v0 s'obtient en additionnant les éléments de la colonne de la matrice de probabilité et les diviser par la somme des éléments de la matrice de probabilité de transition (opération de normalisation). En sachant que la somme des éléments des colonnes doit être égale à zéro.

$$\begin{pmatrix} \frac{0}{10} & \frac{2}{10} & \frac{3}{10} & \frac{5}{10} & \frac{0}{10} \\ \frac{1}{11} & \frac{0}{11} & \frac{4}{11} & \frac{2}{11} & \frac{4}{11} \\ \frac{11}{2} & \frac{11}{4} & \frac{11}{0} & \frac{11}{3} & \frac{11}{3} \\ \frac{12}{3} & \frac{12}{5} & \frac{12}{2} & \frac{12}{0} & \frac{12}{1} \\ \frac{11}{3} & \frac{11}{3} & \frac{11}{3} & \frac{11}{3} & \frac{11}{0} \\ \frac{12}{12} & \frac{12}{12} & \frac{12}{12} & \frac{12}{12} & \frac{12}{12} \end{pmatrix}$$

$\alpha \quad \beta \quad \gamma \quad \phi \quad \lambda$

$\sum \alpha + \sum \beta + \sum \gamma + \sum \phi + \sum \lambda = 1 \rightarrow$ la somme des colonnes doit être égale à 1, raison pour laquelle on devra diviser cette somme par la somme des éléments de la matrice de probabilité de transition.

Ce qui nous donne comme vecteur initial :

$$v_0 = \begin{pmatrix} \frac{1030}{6600} \\ \frac{1634}{6600} \\ \frac{1446}{6600} \\ \frac{1560}{6600} \\ \frac{930}{6600} \end{pmatrix}$$

Nous allons réaliser cette formule itérative jusqu'à convergence. Nous allons définir cette convergence comme étant :

$$|v_1 - v_0| < \textit{gamma}$$

gamma : notre marge d'erreur initialisé à 10^{-6} .

Remarque : Nous prenons la valeur absolue qui est par convention utilisée dans les méthodes numériques pour approximer une solution

Pour ce qui est des itérations de la PowerMethod, elles sont réalisées au point 4 des feuilles scannées.

Conclusion

Nous pouvons donc en conclure qu'avec la combinaison de la PowerMethod et l'utilisation de la matrice/formule de Google nous avons su créer un code/algorithmes qui classe l'ordre d'importance des pages web sur Internet.

Annexe : Résolution du système d'équation linéaire à la main

Université Catholique de
Louvain La neuve

Pani Elio, 6154140
Saley Abdou Djafarou, 81031600
Le 24 décembre 2017

LSINF12.50 : Page 2/10

- ① Trouver la matrice d'adjacence de notre graph. L'indice "i" correspondra aux lignes de la matrice (nœud de départ) et l'indice "j" aux colonnes de la matrice (nœud d'arrivée).

$$A = \begin{matrix} & \begin{matrix} j_1 & j_2 & j_3 & j_4 & j_5 \end{matrix} \\ \begin{matrix} i_1 \\ i_2 \\ i_3 \\ i_4 \\ i_5 \end{matrix} & \begin{bmatrix} 0 & 2 & 3 & 5 & 0 \\ 1 & 0 & 4 & 2 & 4 \\ 2 & 4 & 0 & 3 & 3 \\ 3 & 5 & 2 & 0 & 1 \\ 3 & 3 & 3 & 3 & 0 \end{bmatrix} \end{matrix}$$

la valeur à l'indice A_{ij} correspond au poids de l'arête.

- ② Trouver la matrice de probabilité de transition. je vais utiliser la formule de probabilité conditionnelle du slide 12.3

$$P(\text{page}(k+1) = j \mid \text{page}(k) = i) = \frac{w_{ij}}{w_{i\cdot}} \text{ avec } w_{i\cdot} = \sum_{j=1}^m w_{ij}$$

« la probabilité de passer d'une page i à une page j est égale à la probabilité de passer de la page i à la page j divisé par la probabilité de passer de la page i à n'importe quelle autre page du graph » somme

$$P = \begin{matrix} & \begin{matrix} j_1 & j_2 & j_3 & j_4 & j_5 \end{matrix} \\ \begin{matrix} i_1 \\ i_2 \\ i_3 \\ i_4 \\ i_5 \end{matrix} & \begin{bmatrix} \frac{0}{10} & \frac{2}{10} & \frac{3}{10} & \frac{5}{10} & \frac{0}{10} \\ \frac{1}{11} & \frac{0}{11} & \frac{4}{11} & \frac{2}{11} & \frac{4}{11} \\ \frac{2}{12} & \frac{4}{12} & \frac{0}{12} & \frac{3}{12} & \frac{3}{12} \\ \frac{3}{11} & \frac{5}{11} & \frac{2}{11} & \frac{0}{11} & \frac{1}{11} \\ \frac{3}{12} & \frac{3}{12} & \frac{3}{12} & \frac{3}{12} & \frac{0}{12} \end{bmatrix} \end{matrix}$$

③ Remplissons la matrice de Google avec la formule suivante

$$G_{ij} = \alpha^* P_{ij} + (1-\alpha) * \frac{1}{N} \text{ avec } \alpha = 0,8$$

N : nombre de nœud dans le graph
on obtient :

$$G = \begin{bmatrix} \frac{1}{50} & \frac{9}{50} & \frac{27}{100} & \frac{9}{20} & 0 \\ \frac{9}{110} & \frac{1}{50} & \frac{18}{55} & \frac{9}{55} & \frac{18}{55} \\ \frac{3}{20} & \frac{3}{10} & \frac{1}{50} & \frac{9}{40} & \frac{9}{40} \\ \frac{27}{110} & \frac{9}{22} & \frac{9}{55} & \frac{1}{50} & \frac{9}{110} \\ \frac{9}{40} & \frac{9}{40} & \frac{9}{40} & \frac{9}{40} & \frac{1}{50} \end{bmatrix}$$

"j'ai utilisé le logiciel "matrice calculatrice" en ligne en implémentant la formule " $0,8 * P + (1-0,8) * \frac{1}{5}$ ".

④ nous allons implémenter la formule de la power method :

$$G^* r_0 = r_1$$

Comment trouver le vecteur r_0 ? nous allons normaliser la somme des colonnes de la matrice de probabilité de transition et égaliser l'expression à 1.

$$P = \begin{bmatrix} \frac{0}{10} & \frac{2}{10} & \frac{3}{10} & \frac{5}{10} & \frac{0}{10} \\ \frac{1}{11} & \frac{0}{11} & \frac{4}{11} & \frac{2}{11} & \frac{4}{11} \\ \frac{2}{12} & \frac{4}{12} & \frac{0}{12} & \frac{3}{12} & \frac{3}{12} \\ \frac{3}{11} & \frac{5}{11} & \frac{2}{11} & \frac{0}{11} & \frac{1}{11} \\ \frac{3}{12} & \frac{3}{12} & \frac{3}{12} & \frac{3}{12} & \frac{0}{12} \end{bmatrix} \begin{matrix} i_1 \\ i_2 \\ i_3 \\ i_4 \\ i_5 \end{matrix}$$

$$\alpha = \sum_{\phi=1}^5 \sum_{j=1}^5 \frac{1}{j} \phi + i \lambda$$

α correspond à la somme des éléments de chacune des colonnes.

pour que $\alpha = 1$, il faut que l'ordonnée α par la somme des éléments de la matrice P.

$$\alpha \stackrel{?}{=} 1 \text{ ou } \frac{\alpha}{\sum_{i=1}^5 \sum_{j=1}^5 P_{ij}} = 1$$

ici nous avons à dire que chaque élément a_{ij} du vecteur initial r_0 sera un terme de la somme du membre de gauche de l'expression \Rightarrow

$$\frac{\alpha}{\sum_{i=1}^5 \sum_{j=1}^5 P_{ij}}$$

la somme des éléments de la matrice de probabilité donnée: $\frac{6600}{1320}$

- la somme de la colonne j_1 : $\frac{0}{10} + \frac{1}{11} + \frac{2}{12} + \frac{3}{11} + \frac{3}{12} = \frac{1030}{1320}$

- " " " j_2 : $\frac{2}{10} + 0 + \frac{4}{12} + \frac{5}{11} + \frac{3}{12} = \frac{1634}{1320}$

- " " " j_3 : $\frac{3}{10} + \frac{4}{11} + 0 + \frac{2}{11} + \frac{3}{12} = \frac{1446}{1320}$

- " " " j_4 : $\frac{5}{10} + \frac{2}{11} + \frac{3}{12} + 0 + \frac{3}{12} = \frac{1560}{1320}$

- " " " j_5 : $0 + \frac{4}{11} + \frac{3}{12} + \frac{1}{11} + 0 = \frac{930}{1320}$

la vector r_0 nous donne

$$r_0 = \begin{bmatrix} j_1 \\ j_2 \\ j_3 \\ j_4 \\ j_5 \end{bmatrix} = \begin{bmatrix} \frac{1030}{1320} \\ \frac{1634}{1320} \\ \frac{1446}{1320} \\ \frac{1560}{1320} \\ \frac{930}{1320} \end{bmatrix}$$

mais allons maintenant faire 3 itérations avec la même méthode

On prend $\Rightarrow G^T \cdot r_0 = r_1$ (j'utilise le logiciel en ligne "Matrice calculatrice" pour effectuer les multiplications).

$$\begin{bmatrix} 1 & 9 & 3 & 27 & 9 \\ 50 & 110 & 20 & 110 & 40 \\ 9 & 1 & 3 & 9 & 9 \\ 50 & 50 & 10 & 22 & 40 \\ 27 & 18 & 1 & 9 & 9 \\ 100 & 55 & 50 & 55 & 40 \\ 9 & 9 & 9 & 1 & 9 \\ 20 & 55 & 40 & 50 & 40 \\ 0 & 18 & 9 & 9 & 1 \\ & 55 & 40 & 110 & 50 \end{bmatrix} \cdot \begin{bmatrix} \frac{1030}{1320} \\ \frac{1634}{1320} \\ \frac{1446}{1320} \\ \frac{1560}{1320} \\ \frac{930}{1320} \end{bmatrix} = \begin{bmatrix} 19267 \\ 132000 \\ 1649243 \\ 7260000 \\ 478979 \\ 2420000 \\ 9509 \\ 48400 \\ 6703 \\ 44000 \end{bmatrix}$$

$\Rightarrow G^T \cdot r_0 = r_1$

$$\begin{bmatrix} 1 & 9 & 3 & 27 & 9 \\ 50 & 110 & 20 & 110 & 40 \\ 9 & 1 & 3 & 9 & 9 \\ 50 & 50 & 10 & 22 & 40 \\ 27 & 18 & 1 & 9 & 9 \\ 100 & 55 & 50 & 55 & 40 \\ 9 & 9 & 9 & 1 & 9 \\ 20 & 55 & 40 & 50 & 40 \\ 0 & 18 & 9 & 9 & 1 \\ & 55 & 40 & 110 & 50 \end{bmatrix} \cdot \begin{bmatrix} 19267 \\ 132000 \\ 1649243 \\ 7260000 \\ 478979 \\ 2420000 \\ 9509 \\ 48400 \\ 6703 \\ 44000 \end{bmatrix} = \begin{bmatrix} 427173133 \\ 319440000 \\ 3272252947 \\ 1597200000 \\ 980525551 \\ 5329000000 \\ 14413607 \\ 266200000 \\ 146945341 \\ 106480000 \end{bmatrix}$$

$\Rightarrow G^T \cdot r_1 = r_2$

$$\begin{bmatrix} 1 & 9 & 3 & 27 & 9 \\ 50 & 110 & 20 & 110 & 40 \\ 9 & 1 & 3 & 9 & 3 \\ 50 & 50 & 10 & 22 & 40 \\ 27 & 18 & 1 & 9 & 3 \\ 100 & 53 & 50 & 55 & 40 \\ 9 & 9 & 9 & 1 & 9 \\ 20 & 55 & 40 & 50 & 40 \\ 0 & 18 & 9 & 9 & 1 \\ & 55 & 40 & 110 & 50 \end{bmatrix} \times \begin{bmatrix} 427173133 \\ 3194400000 \\ 3272258347 \\ 1597800000 \\ 380525551 \\ 5324000000 \\ 49413607 \\ 266200000 \\ 14645341 \\ 1064800000 \end{bmatrix} = \begin{bmatrix} 79014152617 \\ 63888000000 \\ 6690630913 \\ 351384000000 \\ 1970854407563 \\ 1171280000000 \\ 139504134897 \\ 1171280000000 \\ 26925780859 \\ 212960000000 \end{bmatrix}$$

$\equiv \quad G^T \quad \times \quad r_2 \quad = \quad r_3$

Le vecteur r_3 est notre vecteur de score PageRank après 3 itérations

Code Python

Méthode pageRankScore

```

1  PageRank.py
2  import numpy as np
3
4  def pageRankScore(mat1, alpha):
5      #we initialize an empty matrix (full of zero) that will be the transition probability matrix
6      transitionMatrix = np.zeros((len(mat1), len(mat1)))
7      #the loop for the transition matrix line's
8      for line in range(0, len(mat1)):
9          #the loop for the transition matrix column's
10         for column in range(0, len(mat1)):
11             #we apply the matrix probability method to get all the value in the transitionMatrix
12             transitionMatrix[line][column] = mat1[line][column] / np.sum(mat1[line])
13
14     #after checking that the matrix was stochastic (the sum of every line =1), let's do the Google matrix
15     #the formula of the Google Matrix : G[i][j] = alpha*transitionMatrix[i][j] + (1-alpha)*(1/N) where N = number of pages/nodes
16     googleMatrix = np.zeros((len(transitionMatrix), len(transitionMatrix)))
17     #we apply the formula above
18     for line in range(0, len(googleMatrix)):
19         for column in range(0, len(googleMatrix)):
20             googleMatrix[line][column] = alpha*transitionMatrix[line][column] + (1-alpha)*(1.0/len(googleMatrix))
21
22     #let's now do the powerMethod in order to find the eigen vector of the google matrix which is simply the pageRankScore of each page
23     #first, we create the vector that we need to multiply with the google matrix, this vector will be the final pageRank solution
24     vectorPageRank = np.zeros((len(mat1), 1))
25     #we make the addition of every column in the adjacency matrix and we store the result in a list
26     listOfColumnSum = mat1.sum(axis = 0)
27     #we fill the initial vector of PageRank with the correct value stored in the list listOfColumnSum
28     for fillin in range(0, len(listOfColumnSum)):
29         vectorPageRank[fillin][0] = listOfColumnSum[fillin] / np.sum(mat1)
30
31     #Before computing the power method, we initialise an error margin which will indicates us that the method found a convergence point
32     gamma = 0.000001
33     #we create a boolean to go out of the loop as soon as we get a good result, I prefer use a boolean instead of a break statement
34     finish = False
35     while(not(finish)):
36         #we store the value of the PageRank vector at the next iteration we take the transpose of the google
37         nextVectorPageRank = np.dot(np.transpose(googleMatrix), vectorPageRank)
38         #we compare the vector at the previous iteration and at the next iteration to see whether the solution is stabilized, the method .all() compare every elements of each ve
39         if (abs(nextVectorPageRank - vectorPageRank < gamma)).all():
40             #we change the boolean of the while loop
41             finish = True
42             #we return the last iteration with the highest accuracy
43             print(nextVectorPageRank)
44         #if the error degree is still to big, we go for an other iteration
45         vectorPageRank = nextVectorPageRank

```

Méthode main

```
45 def main():
46     # we open the file with it's location, in this case, it's in the same folder as the code , no need to print the path
47     fichier = open("matrice-adjacence.csv")
48     #we are going to read the first line of the document
49     contenu = fichier.readline()
50     # the dimension of the matrix according to the previous line (we divide it by 2 because the comas are counted)
51     dimOfMatrix = int(len(contenu)/2)
52     # we initialise a empty matrix (full of zeros)
53     matrice = np.zeros(shape=(dimOfMatrix,dimOfMatrix))
54     # we create a list for a later use which will count 2 by 2 until the end of the line which is read (to avoid the comas)
55     advance = list(range(0,len(contenu),2))
56     # this variable will help us to keep track of the matrix line that we are filling
57     line = 0
58     # that will be our stop condition, whend the variable "contenu" refers to an empty string, that mean that we've reached the end of the file
59     while not(contenu == ""):
60         # we set up the column to the matrix at zero at each iteration to be able to fill the matrix iteratively
61         column = 0
62         for j in advance:
63             # we fill the matrix with the elements inside "contenu", j grows 2 by 2 because we don't need comas in our matrix!
64             matrice[line][column]=float(contenu[j])
65             #we increment the column to go to the next one
66             column +=1
67         line +=1
68         # that's the statement that will continue or not the loop
69         contenu = fichier.readline()
70     # to avoid possible future problems, we have to close the file
71     fichier.close()
72     pageRankScore(matrice, alpha=0.9)
73
74 main()
```