



RAPPORT PROJET FINAL : GRAPHER (MODELÈS ET ALGORITHMES)

EA MAP572 2019-2020

29 novembre 2019

Achraf AZIZE

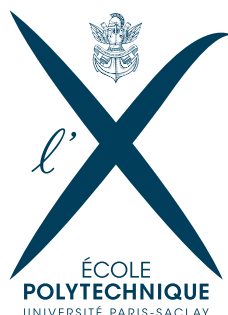


TABLE DES MATIÈRES

1	Abstract	3
2	Graphes par attachement préférentiel	4
2.1	Le Modèle	4
2.2	Simulation	4
2.3	Loi des degrés des sommets	5
2.4	Rich-gets-richer Phenomenon	7
3	Visualisation du graphe	8
3.1	Visualisation aléatoire	8
3.2	Attachement δ -préférentiel	11
4	Spectral Clustering	12
4.1	L'algorithme	12
4.2	Stochastic Block Model	14
5	PageRank	16
5.1	L'algorithme	16
5.2	L'implémentation	16
5.3	Tricher	17

1

ABSTRACT

Ce document constitue un rapport sur les travaux réalisés pendant les trois dernières séances du EA MAP572, dans le cadre du projet final, concernant différents algorithmes et modèles sur les graphes.

Toutes les quatre sections du TP seront abordés, à savoir :

- Les graphes à attachement préférentiel : modèle, simulation et
- Coder un programme qui va calculer les prix pour les contrats actuels et renvoyer un fichier CSV identique à celui fournit par iVector, en moins de temps

Ce document est totalement indépendant, avec tous les codes et scripts utiles. Cependant, un notebook jupyter est joint à ce rapport, pour pouvoir executer directement les tests, et tester la performance.

2

GRAPHS PAR ATTACHEMENT PRÉFÉRENTIEL

Nous allons commencer par le graphe à attachement préférentiel, introduit par A.-L. Barabási et R. Albert pour modéliser de façon dynamique la création de grands réseaux d'interactions.

2.1 LE MODÈLE

Soit $n \geq 1$ un entier, nous allons construire de façon dynamique un graphe aléatoire G_n à n sommets $\{1, \dots, n\}$ et n arêtes :

1. On part du graphe G_1 qui contient un unique sommet $\{1\}$, et dont l'unique arête est une boucle $1 \rightarrow 1$;
2. Connaissant G_k , on rajoute un sommet $k + 1$ et une nouvelle arête $k + 1 \rightarrow v_{k+1}$, où $v_{k+1} \in \{1, 2, \dots, k\}$ est aléatoire, indépendant de G_k , et tiré de la façon suivante :

$$\mathbb{P}(v_{k+1} = i) = \frac{\text{degre}(i)}{\sum_j \text{degre}(j)} = \frac{\text{degre}(i)}{2k-1}$$

où **degre(i)** est le nombre d'arêtes **distinctes** qui touchent le sommet i .

2.2 SIMULATION

Voici le script Python pour simuler la matrice d'adjacence d'un graphe à attachement préférentiel, avec **Adjacence(i,j)**=1 si i et j sont voisins dans le graph :

```
1 import numpy as np
2
3 def Gr(n):
4     if n==1:         # cas n=1
5         return [[1]]
6     deg=[0,1]        # initialisation de la liste des degres
7     A=np.zeros((n,n))
8     A[0][0]=1        # connecter 1 avec 1
9     for i in range(1,n): # pour chaque element
10        # on calcule la liste des probabilites cumules
11        p=np.array(np.cumsum(deg))/np.sum(deg)
12        # on tire un nombre alatoire
13        s=rd.random()
```

```

14         # on regarde sa position j dans p
15         j=-1
16         for k in range(len(p)-1):
17             if s>=p[k] and s<p[k+1]:
18                 j=k
19         # on relie i avec j
20         A[i][j]=1
21         A[j][i]=1
22         # on actualise deg
23         deg[j+1]+=1
24         deg.append(1)
25     # on retourne la matrice d'adjacence ainsi que la liste des degres
26     return A,deg[1:]

```

Listing 1 – Simulation de G_n

Quelques Exemples :

$$n = 2 : \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$$

$$n = 5 : \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

Remarques :

Deux petites remarques concernant le script ci-dessus :

- Le choix d'une méthode itérative, qui est justifié par une performance supérieure en gestion de mémoire par rapport à la méthode récursive, qui prend beaucoup plus de temps pour n assez grand.
- Le choix de retourner la liste des degrés, qui permet directement de l'exploiter dans ce qui suit (la distribution de la loi des degrés)

2.3 LOI DES DEGRÉS DES SOMMETS

Pour n et k assez grand (avec $k < n$), on peut montrer que :

$$\mathbb{P}(\text{degre}(s) = k) \approx ck^{-\alpha}$$

Pour une simulation de la loi des degrés des sommets et pour une bonne illustration, on peut utiliser la fonction `seaborn.displot`, qu'on peut installer avec la commande `conda install seaborn`.

Le script de la simulation est le suivant :

```

1 import seaborn as sbn
2
3 def Degre_distribution(n,n_sim):
4     L=np.zeros(n)
5     for i in range(n_sim): # on simule n_sim fois,
6         _,A=Gr(n)          # pour n tres grand une simulation suffit
7         L+=np.array(A)/n_sim # A repr sente la liste des degr s
8     sbn.distplot(L)
9     plt.show()
10    return L
11
12 L=Degre_distribution(1000,10)

```

Listing 2 – Simulation de la loi des degres des sommets

Pour $n = 1000$ et $n_{sim} = 10$ On retrouve le résultat suivant, où c'est assez clair que c'est loi de puissance.

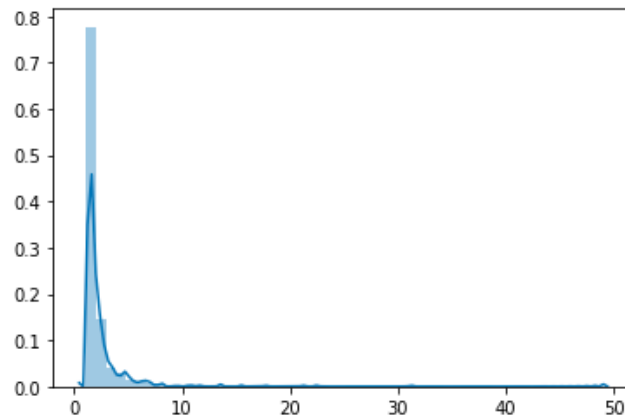


FIGURE 1 – Illustration de la loi de puissance

Reste a savoir comment approximer α .

Pour cela, trois méthodes sont envisageables :

- on pourrait représenter la loi de probabilité dans un repère log-log, dans ce cas, la courbe serait une droite de pente $-\alpha$. Le probleme avec cette méthode c'est que meme pour des n assez grand, y a des degrés qui n'apparaissent pas dans la liste, et donc le log n'est pas défini
- Pour remédier à cela, on peut stocker dans une liste, le log du nombre des apparences non nul, et en utilisant la fonction `optimize.leastsq` pour trouver la meilleur droite qui pourra passer par ces points [1]
- ou plus simplement, utiliser la fonction `powelow.Fit` directement sur la liste des degrés

Une implémentation simple de la dernière méthode donne :

```
1 import powerlaw
2 results = powerlaw.Fit(L)
3 print(results.power_law.alpha)
```

Listing 3 – Evaluation de α

On trouve : $\alpha \approx 3$.

Ce qui est en accord avec les résultats trouvés par Barabási et R.Albert.

Une généralisation :

On peut même généraliser ceci, en prenant :

$$\mathbb{P}(v_{k+1} = i) = \beta \frac{1}{k} + (1 - \beta) \frac{\text{degre}(i)}{2k-1}, \beta \in [0, 1] \setminus \{1\}$$

Et dans ce cas, la loi des degrés serait approximativement, une loi de puissance de parametre $\frac{3-\beta}{1-\beta} \cdot [2]$

Ceci peut donc être une manière concevable pour simuler des lois de puissance.

D'autre part, ceci est aussi une particularité des graphes à attachement préférentiel, qui expliquent l'apparition de la loi de puissance observé aussi dans des modèles dans la nature (réseau sociaux, internet).

2.4 RICH-GETS-RICHER PHENOMENON

On a :

$$E(X_{k+1}|X_k) = p_k(X_{k+1} + 1) + (1 - p_k)X_k$$

avec :

$$p_k = \frac{X_k}{2k-1}$$

car soit on branche $k+1$ avec 1, et donc le nombre de degré devient $X_k + 1$, ou non et le nombre reste le meme.

En passant à l'esperance, on trouve :

$$E(X_{k+1}) = E(X_K) \frac{2k}{2k-1}$$

On voit bien que la moyenne des sommets de 1 ne fait qu'augmenter (strictement croissante), et assez supérieur à la moyenne.

Ainsi, 1 qui commence déjà "riche", le devient de plus en plus.

C'est le Rich-gets-richer Phenomenon.

3

VISUALISATION DU GRAPHE

Dans cette section, il est question de représenter les graphes de manière harmonieuse. Une première approche simple et basique serait de les représenter d'une manière aléatoire.

3.1 VISUALISATION ALÉATOIRE

- On tire au sort n points uniformes indépendants $(X_1, Y_1), \dots, (X_n, Y_n)$ dans le carré $[0, 1]^2$
- Si i, j sont voisins dans G_n , tracer un segment entre (X_i, Y_i) et (X_j, Y_j)

Le script de cette visualisation est le suivant :

```

1 def Graph(A,X,Y):
2     n=len(A)                # A la matrice d'adjacence
3     plt.plot(X,Y,'ob')     # X,Y la liste des absices et ordon es resp
4     for i in range(n):
5         for j in range(i,n):
6             if A[i][j]==1:
7                 plt.plot((X[i],X[j]),(Y[i],Y[j]),'r')
8     plt.show()
9
10 n=50
11 Graph(Gr(n)[0],rd.random(n),rd.random(n))

```

Listing 4 – Evaluation de alpha

Le résultat est une représentation chaotique et désastreuse.

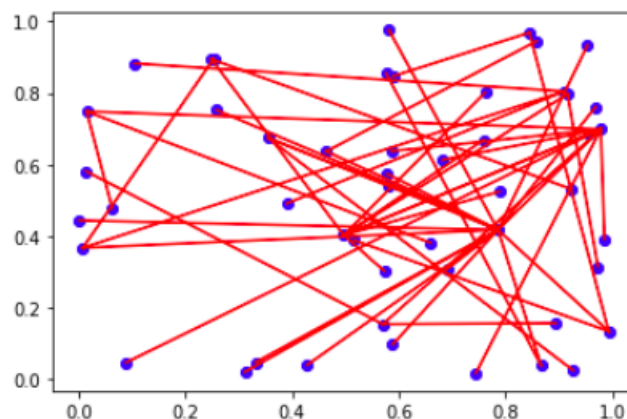


FIGURE 2 – Représentation aléatoire du graphe G_{50}

Pour remédier à cela, on définit l'énergie suivante :

$$E = \sum_{i,j} \frac{\left(\frac{1}{\sqrt{2}}\|M_i - M_j\| - D_{i,j}^*\right)^2}{(D_{i,j}^*)^2}$$

où

$$D_{i,j}^* = \frac{Distance(i,j)}{\max_{i,j} Distance(i,j)}$$

On cherche la configuration des M_i qui va minimiser cette énergie en utilisant l'algorithme "gradient descent". Mais avant, on introduit les fonctions intermédiaires suivantes :

```

1 # cette fonction retourne la matrice D_etoile defini ci-dessus
2 def distance_etoile(G):          #prend en parametre la matrice d'adjacence
3     n=len(G)
4     A=np.zeros((n,n))
5     A[0][1],A[1][0]=1,1          #initialisation
6     for i in range(2,n):
7         for j in range(i):
8             if G[i][j]==1:        #pour chaque iteration, on cherche l'element
                branch
9                 for k in range(i):
10                     A[k][i],A[i][k]=A[k][j]+1,A[k][j]+1  #on actualise la
distance
11                     break
12     return A/np.max(A)
13 #cette fonction retourne la distance entre Mi et Mj
14 def dist_sqr2(X,Y,i,j):
15     return np.sqrt(((X[i]-X[j])**2+(Y[i]-Y[j])**2)/2)
16 def E(G,X,Y,D):
17     n=len(G)
18     S=0
19     for i in range(n):
20         for j in range(i):
21             d=dist_sqr2(X,Y,i,j)
22             S+=((d-D[i][j])**2)/(D[i][j])**2
23     return S
24 #cette fonction retourne le gradient de E
25 def derive_E(G,X,Y,D):
26     n=len(G)
27     X_der=np.zeros(n)
28     Y_der=np.zeros(n)
29     for i in range(n):
30         S,T=0,0
31         for j in range(n):
32             if j!=i:
33                 d=dist_sqr2(X,Y,i,j)
34                 S+=((d-D[i][j])*(X[i]-X[j]))/(d*(D[i][j])**2)
35                 T+=((d-D[i][j])*(Y[i]-Y[j]))/(d*(D[i][j])**2)
36         X_der[i]=S
37         Y_der[i]=T
38     return X_der,Y_der

```

```

39 #cette fonction calcule la variation totale entre deux configurations des Mi
40 def var(X,Y,M,N):
41     return np.sum([np.sqrt((X[i]-M[i])**2+(Y[i]-N[i])**2) for i in range(
        len(X))])

```

Listing 5 – Fonctions intermédiaires

Le script du "gradient descent" est le suivant :

```

1 #epsilon repr sente la limite de convergence, et alpha le pas de la
  descente
2 def Gradient_descent(G,epsilon,alpha):
3     n=len(G)
4     X=rd.random(n)
5     Y=rd.random(n)
6     D=distance_etoile(G)
7     v=epsilon+1
8     while(v>epsilon):
9         X_der,Y_der=derive_E(G,X,Y,D)
10        v=var(X,Y,X-alpha*X_der,Y-alpha*Y_der)
11        X,Y=X-alpha*X_der,Y-alpha*Y_der
12    return X,Y
13
14 X,Y=Gradient_descent(G,0.01,0.001)
15 Graph(G,X,Y)

```

Listing 6 – Gradient descent

Pour le même graphe que tout à l'heure, on trouve cette fois une représentation plus harmonieuse avec $\epsilon=0.01$ et $\alpha=0.001$: Le choix de ϵ contrôle la convergence du modèle.

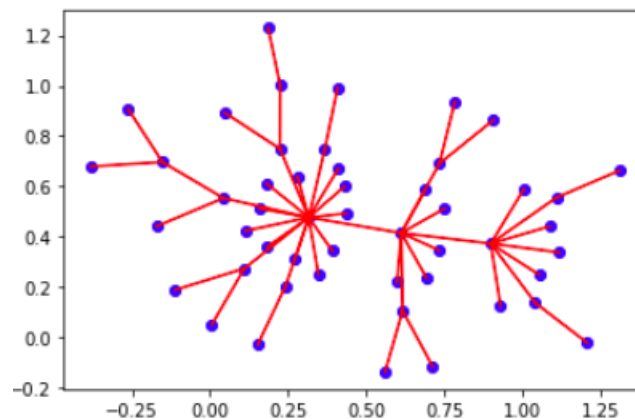


FIGURE 3 – Représentation harmonieuse du graphe G_{50}

Le choix du pas est assez cruciale. Pour un assez grand, le modèle diverge rapidement. D'autre part, on peut aussi représenter ce graphe en 3D, il suffit d'ajouter Z aux fonctions précédente.

3.2 ATTACHEMENT δ -PRÉFÉRENTIEL

On peut généraliser le modèle du graphe à attachement préférentiel de la façon suivante. Soit $\delta > -1$ un paramètre fixé, on remplace l'équation précédente par

$$\mathbb{P}(v_{k+1} = i) = \frac{\text{degre}(i) + \delta}{\sum_j (\text{degre}(j) + \delta)} = \frac{\text{degre}(i) + \delta}{2k - 1 + k\delta}$$

Pour $\delta = 0$, on retrouve le modèle initial. Le script de simulation est le suivant :

```
1 def Gr_delta(n,d):
2     if n==1:
3         return np.array([[1]])
4     deg=[0,1+d]
5     A=np.zeros((n,n))
6     A[0][0]=1
7     for i in range(1,n):
8         p=np.array(np.cumsum(deg)/np.sum(deg))
9         a=rd.random()
10        j=-1
11        for k in range(len(p)-1):
12            if a >=p[k] and a <p[k+1]:
13                j=k
14        deg[j+1]+=1
15        deg.append(1+d)
16        A[i][j]=1
17        A[j][i]=1
18    return A
```

Listing 7 – Gradient descent

Quand δ tend vers l'infini, on choisit v_{k+1} de manière uniforme dans $\{1, \dots, k\}$

Ainsi, δ permet de contrôler la préférence par rapport aux sommets.

Pour $\delta = 0$, on a vu qu'il y a une tendance de "Rich gets richer".

Pour δ plus grand, le graph est plus filaire et homogène (moins de concentration de richesse).

Pour $n = 50$ et $n = 10000$, on trouve :

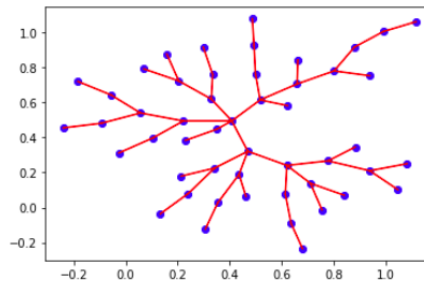


FIGURE 4 – Représentation du graphe δ préférentiel G_{50}

4

SPECTRAL CLUSTERING

4.1 L'ALGORITHME

Soit k le nombre de clusters, et n la taille du graphe.

- On calcule le laplacien $L=D-A$, avec D une matrice diagonale tel que $D_{i,i} = \text{degre}(i)$ et A la matrice d'adjacence
- On calcule les k premiers vecteurs propres de L
- Soit $U \in R^{n \times k}$ la matrice qui contient ces vecteurs en colonnes
- Pour $i = 1, \dots, n$ soit $y_i \in R^k$, on "cluster" ces points en k clusters en utilisant l'algorithme k-means

Le script correspédant est le suivant :

```
1 from sklearn.cluster import KMeans
2 def spectral_cluster(G,k):
3     D = np.diag(G.sum(axis=1))
4     L=D-G
5     vals,vecs = np.linalg.eig(L)
6     vecs = vecs[:,np.argsort(vals)]
7     vals = vals[np.argsort(vals)]
8     kmeans = KMeans(n_clusters=k)
9     kmeans.fit(vecs[:,1:k])
10    label = kmeans.labels_
11    colors=["blue","red","green","cyan","magenta","yellow","black","white"]
12    for i in range(len(G)):
13        plt.scatter(X[i],Y[i],c=colors[label[i]])
14    for i in range(len(G)):
15        for j in range(i,len(G)):
16            if G[i][j]==1:
17                plt.plot((X[i],X[j]),(Y[i],Y[j]),'black')
18    plt.xlabel("k={}".format(k))
19    plt.show()
```

Listing 8 – Spectral Clustering

Les sorties de cette algorithme pour différentes valeurs de k

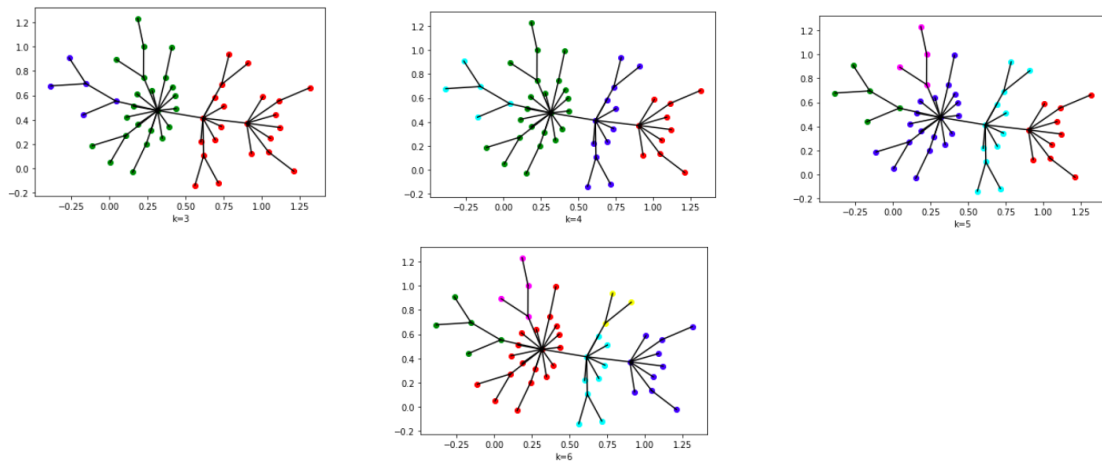


FIGURE 5 – Spectral clustering pour différents k

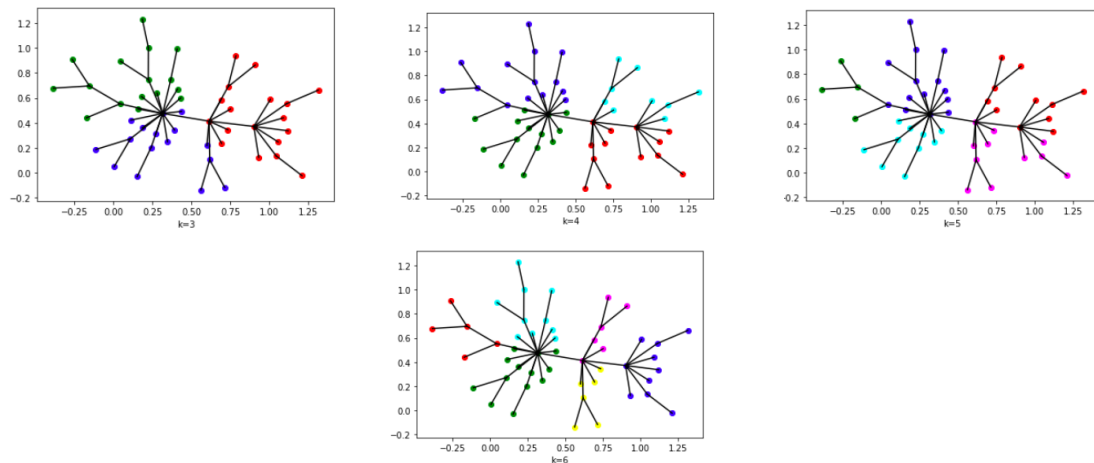


FIGURE 6 – K-means clustering pour différents k

On pourra le comparer avec un K-Means directement appliqué sur les points du graphe, comme dans la Figure 6 ci-dessus.

On voit clairement la particularité du Spectral Clustering.

Le K-means permet de trouver les clusters en minimisant la distance total, sans prendre en compte la structure du graphe.

D'autre part, le Spectral Clustering combine les deux contraintes.

4.2 STOCHASTIC BLOCK MODEL

On va maintenant s'intéresser à implémenter le Spectral Clustering sur un autre modèle de graph.

Soit E_1, \dots, E_k une partition de $\{1, \dots, n\}$ en K clusters, et $(q_{r,s})_{r,s \leq K}$ la matrice de probabilité inter-classe.

Si $i \in E_r$ et $j \in E_s$, on met une arrête entre i et j avec probabilité $q_{r,s}$

Supposons qu'un graphe est déjà simulé avec le modèle SBM, le but est de retrouver les clusters avec le Spectral Clustering.

On le test pour des petites valeurs de n .

```

1 from graspy.simulations import sbm
2 K=2
3 n_SBM = [50, 50]
4 P_SBM = [[0.8, 0.2],
5           [0.2, 0.5]]
6 G_SBM = sbm(n=n_SBM, p=P_SBM)
7 D_SBM = np.diag(G_SBM.sum(axis=1))
8 L_SBM=D_SBM-G_SBM
9 vals_SBM,vecs_SBM = np.linalg.eig(L_SBM)
10 vecs_SBM = vecs_SBM[:,np.argsort(vals_SBM)]
11 vals_SBM = vals_SBM[np.argsort(vals_SBM)]
12 kmeans_SBM = KMeans(n_clusters=K)
13 kmeans_SBM.fit(vecs_SBM[:,1:K])
14 label_SBM = kmeans_SBM.labels_
15 print(label_SBM)

```

Listing 9 – Spectral Clustering sur le SBM

[illegible]

On récupère exactement les deux clusters initiaux, car on l'applique avec un nombre de K déjà connus, et la probabilité de rester dans le même cluster est assez grande (les éléments de la diagonale de la matrice de probabilité).

Si ces éléments sont assez petits, il est plus difficile de récupérer ces clusters.

Il est même possible de retrouver des conditions plus exactes entre ces deux probabilités pour retrouver le modèle.

```
1 MatriceAdjacence=np.loadtxt('StochasticBlockModel.txt')
2 D2 = np.diag(MatriceAdjacence.sum(axis=1))
3 L2=D2-MatriceAdjacence
4 vals2,vecs2 = np.linalg.eig(L2)
5 vecs2 = vecs2[:,np.argsort(vals2)]
6 vals2 = vals2[np.argsort(vals2)]
7 kmeans_ADJ = KMeans(n_clusters=4)
8 kmeans_ADJ.fit(vecs2[:,1:4])
9 label_ADJ = kmeans_ADJ.labels
```

```
10 print(label_ADJ)
```

Listing 10 – Spectral Clustering sur l'exemple

Pour la matrice proposé, on trouve pour $K=4$ la partition suivante : [360 90 60 90]

5

PAGERANK

5.1 L'ALGORITHME

Afin de trouver le score de popularité d'un graphe, on cherche le vecteur propre à gauche de la matrice suivante :

$$P_\varepsilon = (1 - \varepsilon)\tilde{A} + \frac{\varepsilon}{n} \begin{pmatrix} 1 & \dots & 1 \\ & \cdot & \\ 1 & \dots & 1 \end{pmatrix}$$

avec \tilde{A} est la matrice d'adjacence renormalisé.

5.2 L'IMPLÉMENTATION

Le script Python est le suivant :

```
1 from scipy.linalg import eig
2 def P(e,G):
3     return (1-e)*(G/np.sum(G,axis=1)) + e*np.ones((len(G),len(G)))/(len(G))
4 def pagerank(e,G):
5     n=len(G)
6     S=np.zeros((n,n))
7     # on trinagularise
8     for i in range(n):
9         for j in range(i+1):
10             S[i][j]=G[i][j]
11     Pe=P(e,S)
12     values, vecs = eig(Pe, right = False, left = True)
13     vecs = vecs[:,np.argsort(values)]
14     values = values[np.argsort(values)]
15     for i in range(n):
16         if round(values[i].real,2)==1.00 :
17             return vecs[:,i].real/np.sum(vecs[:,i].real)
18     return np.zeros(n)
```

Listing 11 – Pagerank

Remarques :

- Sans triangulariser, l'algorithme ne fournit pas des résultats correctes
- pour epsilon=0, on trouve [1,0,...,0]

- pour $\epsilon=1$, on trouve $[1,1,\dots,1]$ normalisé
- ϵ représente donc le degré d'exploration du graphe
- prendre $\epsilon=0.15$ permet d'avoir un bon critère de popularité
- il est quand même toujours possible d'augmenter son score en créant artificiellement quelques sommets qui pointent vers un sommet fixé

5.3 TRICHER

```
1 def tricher(e,G):
2     n=len(G)
3     S=np.zeros((n+3,n+3))
4     for i in range(n):
5         for j in range(n):
6             S[i][j]=G[i][j]
7     S[n][0],S[n+1][0],S[n+2][0]=1,1,1
8     return pagerank(e,S)[0]-pagerank(e,G)[0]
9 tricher(0.15,F)
```

Listing 12 – Augmenter le score PageRank

Dans cet exemple, on ajoute manuellement 3 sommets avec 1, et on remarque que son score augmente, ce qui est normal.

RÉFÉRENCES

- [1] Fitting data in python. <https://scipy-cookbook.readthedocs.io/items/FittingData.html>.
- [2] Power-law random graphs. <https://www.di.ens.fr/~lelarge/X15/cours/slides8.pdf>.