

Algorithmique des Réseaux Sociaux

21 Octobre

Exercice 1 – PageRank et temps de mixage

On définit la distance en variation totale entre deux mesures de probabilité μ et ν sur Ω par :

$$\|\mu - \nu\|_{\text{TV}} = \max_{A \subset \Omega} |\mu(A) - \nu(A)|.$$

1. Montrer que

$$\frac{1}{2} \sum_{x \in \Omega} |\mu(x) - \nu(x)| = \sum_{\mu(x) \geq \nu(x)} (\mu(x) - \nu(x)).$$

2. Montrer que pour $B = \{x, \mu(x) \geq \nu(x)\}$ et tout $A \subset \Omega$, on a :

$$\mu(A) - \nu(A) \leq \mu(B) - \nu(B)$$

3. En déduire que

$$\|\mu - \nu\|_{\text{TV}} = \frac{1}{2} \sum_{x \in \Omega} |\mu(x) - \nu(x)|.$$

Un couplage entre deux mesures de probabilité μ et ν est une paire de variables aléatoires (X, Y) définies sur un même espace de probabilité telle que la distribution marginale de X est μ , i.e. $\mathbb{P}(X = x) = \mu(x)$ et la marginale de Y est ν , i.e. $\mathbb{P}(Y = y) = \nu(y)$.

4. Montrer que

$$\|\mu - \nu\|_{\text{TV}} \leq \inf\{\mathbb{P}(X \neq Y) : (X, Y) \text{ est un couplage de } \mu \text{ et } \nu.\}$$

5. Montrer que

$$\sum_{x \in \Omega} \mu(x) \wedge \nu(x) = 1 - \|\mu - \nu\|_{\text{TV}} \in [0, 1].$$

6. En déduire un couplage (X, Y) tel que $\|\mu - \nu\|_{\text{TV}} = \mathbb{P}(X \neq Y)$.

Une chaîne de Markov à espace d'états Ω et de matrice de transition P est une suite de variables aléatoires (X_0, X_1, \dots) telles que

$$\mathbb{P}(X_{t+1} = y | X_t = x, X_{t-1} = x_{t-1}, \dots, X_0 = x_0) = \mathbb{P}(X_{t+1} = y | X_t = x) = P(x, y).$$

Une chaîne est dite irréductible si pour tous $x, y \in \Omega$, il existe t tel que $P^t(x, y) > 0$. Soit $T(x) = \{t \geq 1, P^t(x, x) > 0\}$. La période de x est le pgcd de $T(x)$. Une chaîne est dite apériodique si tous les états ont pour période 1.

7. Montrer que si P est irréductible alors tous les états ont même période.

8. Montrer que si P est irréductible apériodique alors P est primitive.

9. Montrer que dans ce cas, il existe une unique distribution stationnaire sur Ω , satisfaisant $\pi = \pi P$.

On note $P^t(x, \cdot)$ pour la loi de la chaîne de Markov au temps t commencée en $X_0 = x \in \Omega$ et μP^t pour le cas où la loi de X_0 est μ . On définit

$$\begin{aligned} d(t) &= \max_{x \in \Omega} \|P^t(x, \cdot) - \pi\|_{\text{TV}}, \\ \bar{d}(t) &= \max_{x, y \in \Omega} \|P^t(x, \cdot) - P^t(y, \cdot)\|_{\text{TV}}. \end{aligned}$$

10. Montrer que

$$d(t) \leq \bar{d}(t) \leq 2d(t).$$

11. Montrer que

$$\begin{aligned} d(t) &= \sup_{\mu} \|\mu P^t - \pi\|_{\text{TV}}, \\ \bar{d}(t) &= \sup_{\mu, \nu} \|\mu P^t - \nu P^t\|_{\text{TV}}, \end{aligned}$$

où μ et ν sont des distributions sur Ω .

12. Montrer que

$$\|\mu P - \nu P\|_{\text{TV}} \leq \|\mu - \nu\|_{\text{TV}}.$$

En déduire que $d(t)$ et $\bar{d}(t)$ sont décroissantes en t .

13. Pour un couplage (X_s, Y_s) de $P^s(x, \cdot)$ et de $P^s(y, \cdot)$, montrer que

$$\|P^{t+s}(x, \cdot) - P^{t+s}(y, \cdot)\|_{\text{TV}} = \frac{1}{2} \sum_w |\mathbb{E}[P^t(X_s, w) - P^t(Y_s, w)]|.$$

14. En déduire que \bar{d} est sous-multiplicative : $\bar{d}(t+s) \leq \bar{d}(t)\bar{d}(s)$.

On définit le temps de mixage par

$$t_{\text{mix}}(\epsilon) = \min\{t, d(t) \leq \epsilon\} \text{ et } t_{\text{mix}} = t_{\text{mix}}(1/4).$$

15. Montrer que $d(kt_{\text{mix}}) \leq 2^{-k}$ et $t_{\text{mix}}(\epsilon) \leq \lceil \log_2 \epsilon^{-1} \rceil t_{\text{mix}}$.

On définit un couplage de chaînes de Markov avec matrice de transition P comme le processus $(X_t, Y_t)_{t=0}^{\infty}$ ayant la propriété que chacun des processus (X_t) et (Y_t) est une chaîne de Markov de matrice de transition P (ces processus peuvent avoir des points de départs différents). Tout couplage peut être modifié de telle sorte que : si $X_s = Y_s$ alors $X_t = Y_t$ pour tout $t \geq s$. On définit alors le temps de couplage par

$$\tau_{\text{couple}} = \min\{t, X_t = Y_t\}.$$

16. Comme exemple, considérer la marche aléatoire sur $\{0, 1, \dots, n\}$ qui monte ou descend avec probabilité $1/2$ et qui aux bords, reste immobile avec probabilité $1/2$. En construisant un couplage, montrer que $P^t(x, n) \leq P^t(y, n)$ dès que $x \leq y$.

17. Montrer que

$$\|P^t(x, \cdot) - P^t(y, \cdot)\|_{\text{TV}} \leq \mathbb{P}_{x, y}(\tau_{\text{couple}} > t).$$

18. Considérer la chaîne de Markov associée au calcul de PageRank (qui choisit avec probabilité $1 - \alpha$ un sommet uniformément au hasard à chaque étape). Montrer que son temps de mixage satisfait $t_{\text{mix}} \leq \frac{4}{1-\alpha}$. Qu'en déduire sur la vitesse de convergence de l'algorithme ?

Exercice 2 – PageRank : algèbre linéaire

Comme dans le cours, on note M la matrice carrée telle que $M_{ij} = \frac{1}{d^{+(i)}}$ si $i \rightarrow j$ et zero sinon. Si i est une feuille, i.e. n'a pas de lien sortant alors la ligne i de la matrice M est nulle. On définit \overline{M} la matrice carrée où chacune de ces lignes nulles est remplacée par une ligne $(\frac{1}{n}, \dots, \frac{1}{n})$.

1. Montrer que 1 est valeur propre de \overline{M} et toute valeur propre (complexe) λ de \overline{M} vérifie $|\lambda| \leq 1$.

On définit alors $\tilde{M} = \alpha \overline{M} + (1 - \alpha) \mathbf{1} \mathbf{p}^T$ où \mathbf{p} est un vecteur de probabilité. Remarquer que le cas $\mathbf{p} = \frac{1}{n} \mathbf{1}$ est le cas vu en cours.

2. Montrer qu'il existe une matrice non-singulière $P = (\mathbf{1} X)$ telle que :

$$P^{-1} \tilde{M} P = \begin{pmatrix} 1 & A \\ 0 & B \end{pmatrix}.$$

3. En utilisant P , montrer que les valeurs propres de \tilde{M} sont $\{1, \alpha \lambda_2, \dots, \alpha \lambda_n\}$ où $\{1, \lambda_2, \dots, \lambda_n\}$ sont les valeurs propres de \overline{M} .

4. Conclure sur la vitesse de convergence de l'algorithme PageRank.

On définit le vecteur \mathbf{a} par $a_i = 1$ si la ligne i de M correspond à une feuille, et 0 sinon. Montrer que l'algorithme PageRank calcule π qui satisfait $\pi^T (I - \alpha M) = \mathbf{p}^T$ avec $\pi^T \mathbf{1} = 1$.

5. Montrer que $I - \alpha M$ est inversible pour $\alpha < 1$.

6. En écrivant

$$M = \begin{pmatrix} M_{11} & M_{12} \\ 0 & 0 \end{pmatrix},$$

montrer que $I - \alpha M_{11}$ est inversible puis que l'algorithme suivant calcule le vecteur π :

(a) résoudre : $\pi_1^T (I - \alpha M_{11}) = \mathbf{p}_1^T$.

(b) calculer : $\pi_2^T = \alpha \pi_1^T P_{12} + v_2^T$.

(c) normaliser $\pi = \frac{1}{\|\pi_1\|_1 + \|\pi_2\|_1} (\pi_1, \pi_2)$.

Quel est l'intérêt de cet algorithme ?