

TD 05 – Inégalité de Chernoff et partiel de l'an dernier (corrigé)

Exercice 1.*BucketSort*

Le tri par seaux est un algorithme de tri très simple dont la complexité moyenne est linéaire. Les hypothèses sont les suivantes : nous avons $n = 2^m$ nombres à trier, tirés uniformément et indépendamment sur $\{0, \dots, 2^k - 1\}$ avec $k \geq m$ (k est supposé connu). L'algorithme procède en deux étapes : d'abord, il effectue un pré-tri en jetant selon certaines règles les n éléments dans n seaux ; ensuite, il appelle un algorithme de tri simple (en temps quadratique, tel que tri par insertion ou tri par sélection) dans chaque seau. Enfin, il concatène dans l'ordre les listes triées obtenues dans chaque seau. Pour que l'algorithme soit exact, il faut que le pré-tri soit fait de telle sorte que tous les éléments du seau i soient inférieurs à tous les éléments du seau j , pour $i < j$.

- Donner une façon très simple de faire le pré-tri tout en respectant la condition énoncée ci-dessus. On veut que le choix du seau pour un élément x soit effectué en temps constant (on suppose ici que les opérations arithmétiques peuvent être effectuées en temps constant).

☞ On suppose que les seaux sont numérotés de 0 à $n-1$. L'élément x peut être vu comme un nombre en binaire écrit sur k bits. On regarde les m bits de poids forts (ce qui revient à diviser par 2^{k-m}), cela nous donne un nombre s_x entre 0 et $n-1 \rightarrow$ on met x dans le seau numéro s_x . Ainsi, on est sûr que $x < y$ si $s_x < s_y$, donc la condition d'exactitude du pré-tri est respectée.

- Soit X_i la v.a. comptant le nombre d'éléments dans le seau i après le pré-tri. Quelle loi suit X_i ?

☞ X_i suit une loi binomiale de paramètres $(n, 1/n)$, car les n éléments d'entrée sont choisis indépendamment et uniformément.

- Prouver que la complexité en moyenne est $\mathcal{O}(n)$.

☞ Le pré-tri coûte un temps constant pour chaque élément, donc un temps linéaire en n au total. Ensuite, le tri du seau numéro i coûte $c(X_i)^2$ pour une certaine constante c . Or, comme X_i suit une loi binomiale $(n, 1/n)$, on a (avec $p = 1/n$, $\mathbb{E}[X_i] = np$, $\text{Var}[X_i] = np(1-p)$) :

$$\mathbb{E}[X_i^2] = \mathbb{E}[X_i]^2 + \text{Var}[X_i] = (np)^2 + np(1-p) = np(np+1-p) = 1 \cdot (1+1-\frac{1}{n}) = 2 - \frac{1}{n} < 2.$$

Donc l'espérance du temps passé dans la deuxième étape est au plus

$$\mathbb{E}\left[\sum_{i=0}^{n-1} c(X_i)^2\right] = c \sum_{i=0}^{n-1} \mathbb{E}[X_i^2] \leq 2cn.$$

L'espérance du temps d'exécution totale est donc linéaire en n .

Exercice 2.*Sondage*

On veut faire un sondage d'opinion pour estimer la proportion p de la population qui est en accord avec le Président. Supposons que l'on interroge n personnes choisies uniformément et indépendamment au hasard, et que chacune d'elle réponde par "Oui, je suis d'accord" ou "Non, je ne suis pas d'accord". Étant donné $\theta > 0$ et $0 < \delta < 1$, on souhaite trouver une estimation \bar{X} de p telle que

$$\mathbf{P}\{|\bar{X} - p| \leq \theta\} > 1 - \delta.$$

Par exemple pour $1 - \delta = 0.95$, on pourra ainsi dire que le sondage a une précision de θ à 95%.

- Que choisir comme estimation \bar{X} de p ?

☞ Soit $X_i = 1$ si la i ème personne interrogée est d'accord, et 0 zéro. On pose ensuite $X = \sum_{i=1}^n X_i$ et $\bar{X} = 1/n \cdot X$. Alors $\mathbb{E}[X] = np$ et $\mathbb{E}[\bar{X}] = p$.

- Combien de personnes doit-on interroger pour que l'estimation \bar{X} vérifie nos conditions ? Autrement dit, donner une borne inférieure sur n en termes de θ et δ . On remarquera que cette borne ne dépend pas de la taille de la population totale.

☞ On utilise la borne de Chernoff 'two-sided' sur X :

$$\mathbf{P}\{|\bar{X} - p| \geq \varepsilon p\} = \mathbf{P}\{|X - pn| \geq \varepsilon pn\} \leq 2 \exp\left(-\frac{\varepsilon^2}{2 + \varepsilon} \cdot pn\right).$$

Ensuite on pose $\varepsilon = \theta/p$ pour avoir $\varepsilon p = \theta$, et en ré-injectant :

$$\mathbb{P}\{|\bar{X} - p| \geq \varepsilon p\} \leq 2 \exp\left(-\frac{\theta^2/p^2}{2+\theta/p} \cdot pn\right) = 2 \exp\left(-\frac{\theta^2}{2p+\theta} \cdot n\right).$$

Essayons maintenant de borner ceci par δ . On a :

$$\frac{\theta^2}{2p+\theta} \geq \frac{\theta^2}{2+\delta}$$

et donc

$$\mathbb{P}\{|\bar{X} - p| \geq \theta\} \leq 2 \exp\left(-\frac{\theta^2}{2+\theta} \cdot n\right).$$

Enfin :

$$\begin{aligned} \delta \geq 2 \exp\left(-\frac{\theta^2}{2+\theta} \cdot n\right) &\Leftrightarrow \exp\left(\frac{\theta^2}{2+\theta} \cdot n\right) \geq \frac{2}{\delta} \\ &\Leftrightarrow \frac{\theta^2}{2+\theta} n \geq \ln \frac{2}{\delta} \\ &\Leftrightarrow n \geq \frac{2+\theta}{\theta^2} \ln \frac{2}{\delta} \end{aligned}$$

3. Calculer la valeur de n obtenue grâce à votre borne pour les paramètres $\theta = 0.2$ et $1 - \delta = 95\%$.

☞ On obtient $n \geq 203$.

Exercice 3.

Interrupteurs

Partie I :

1. Montrer qu'il existe une constante $\gamma > 0$ rendant l'énoncé suivant vrai : si une v.a. positive X vérifie $\mathbb{E}[X] = 1$ et $\mathbb{E}[X^2] \leq 3$, alors $\mathbb{P}(X \geq 1/4) \geq \gamma$.

☞ On écrit

$$1 = \mathbb{E}[X] = \mathbb{E}[X \mathbf{1}_{X < 1/4}] + \mathbb{E}[X \mathbf{1}_{X \geq 1/4}] \leq \frac{1}{4} + \mathbb{E}[X \mathbf{1}_{X \geq 1/4}].$$

Par l'inégalité de Cauchy-Schwarz, $\mathbb{E}[X \mathbf{1}_{X \geq 1/4}] \leq \sqrt{\mathbb{E}[X^2] \mathbb{P}(X \geq 1/4)} \leq \sqrt{3} \sqrt{\mathbb{P}(X \geq 1/4)}$. On obtient la minoration voulue pour $\gamma = 3/16$.

2. Soient (X_1, \dots, X_n) des v.a. i.i.d. vérifiant $\mathbb{P}(X_i = 1) = \mathbb{P}(X_i = -1) = \frac{1}{2}$. On pose $Y = \frac{1}{\sqrt{n}}(X_1 + \dots + X_n)$. Calculer $\mathbb{E}[Y^2]$ et $\mathbb{E}[Y^4]$ et en déduire que

$$\mathbb{E}[|X_1 + \dots + X_n|] \geq \frac{\gamma}{2} \sqrt{n}.$$

☞ On a $\mathbb{E}[Y^2] = \frac{1}{n} \cdot \text{Var}[Y] = \frac{1}{n} \cdot \sum_i \text{Var}[X_i] = 1$ (par indépendance). On a ensuite

$$\mathbb{E}[Y^4] = \frac{1}{n^2} \cdot \sum_{i,j,k,l=1}^n \mathbb{E}[X_i X_j X_k X_l].$$

L'indépendance des X_i et le fait que $\mathbb{E}[X_i] = 0$ implique $\mathbb{E}[X_i X_j X_k X_l] = 0$ dès qu'un indice apparaît une unique fois parmi $\{i, j, k, l\}$. Les seuls termes non nuls sont ceux où $i = j = k = l$ ou $i = j \neq k = l$ ou $i = k \neq j = l$ ou $i = l \neq j = k$. On a donc

$$\mathbb{E}[Y^4] = 1/n^2(n + 3n(n-1)) = 3 - 2/n \leq 3.$$

On applique la question précédente à $X = Y^2$, d'où $\mathbb{P}(Y^2 \geq 1/4) = \mathbb{P}(|X_1 + \dots + X_n| \geq \frac{\sqrt{n}}{2}) \geq \gamma$. Enfin,

$$\mathbb{E}[|X_1 + \dots + X_n|] \geq \frac{\sqrt{n}}{2} \mathbb{P}\left(|X_1 + \dots + X_n| \geq \frac{\sqrt{n}}{2}\right) \geq \frac{\gamma \sqrt{n}}{2}.$$

Partie II :

On considère une grille $n \times n$ d'ampoules ainsi que 3 séries d'interrupteurs : des interrupteurs $a = (a_{ij})_{1 \leq i, j \leq n}$ associés à chaque ampoule, des interrupteurs $b = (b_i)_{1 \leq i \leq n}$ associés à chaque ligne et des interrupteurs $c = (c_j)_{1 \leq j \leq n}$ associés à chaque colonne. Chaque interrupteur prend la valeur -1 ou 1 . L'ampoule en position (i, j) est allumée si et seulement si $a_{ij} b_i c_j = 1$. On considère la quantité

$$F(a, b, c) = \sum_{i,j=1}^n a_{ij} b_i c_j$$

qui est le nombre d'ampoules allumées moins le nombre d'ampoules éteintes. Enfin, deux joueurs jouent au jeu suivant : le joueur 1 choisit la position des interrupteurs (a_{ij}) , puis le joueur 2 choisit la position des interrupteurs (b_i) et (c_j) . Le joueur 1 veut minimiser $F(a, b, c)$ et joueur 2 veut le maximiser. On considère donc

$$V(n) = \min_{a \in \{-1,1\}^{n \times n}} \max_{b, c \in \{-1,1\}^n} F(a, b, c).$$

3. Montrer que $V(n) = O(n^{3/2})$ en considérant le cas où le joueur 1 joue au hasard.

☞ Soit $(a_{ij})_{1 \leq i, j \leq n}$ des v.a. i.i.d. de loi uniforme sur $\{-1,1\}$. Quel que soit le choix de b et c , on a

$$\mathbb{P}(F(a, b, c) \geq t) \leq \exp(-t^2/2n^2)$$

par l'inégalité de Chernoff (en effet, $F(a, b, c)$ est la somme de n^2 v.a. de loi uniforme sur $\{-1,1\}$). Par la borne de l'union,

$$\mathbb{P}(\max_{b, c} F(a, b, c) \geq t) \leq 4^n \exp(-t^2/2n^2).$$

Lorsque $t > \sqrt{2n^3 \log 4}$, cette probabilité est < 1 et donc $\mathbb{P}(\max_{b, c} F(a, b, c) < t) > 0$: il existe donc un choix de a tel que $\max_{b, c} F(a, b, c) < t$, d'où $V(n) = O(n^{3/2})$.

4. Le joueur 2 applique la stratégie suivante : il choisit b au hasard, puis ensuite choisit c de façon à allumer le maximum de lampes. Estimer le nombre moyen de lampes allumées par cette stratégie à l'aide de la question I.2 et en déduire que $V(n) = \Omega(n^{3/2})$.

☞ Fixons $a = (a_{ij})$ et choisissons (b_i) i.i.d. de loi uniforme sur $\{-1,1\}$. On a alors

$$\max_c F(a, b, c) = \sum_{i=1}^n \left| \sum_{j=1}^n a_{ij} b_j \right|.$$

En utilisant la linéarité de l'espérance, le fait que $(b_j)_j$ et $(a_{ij} b_j)_j$ ont même loi et la question I.2, il vient

$$\mathbb{E} \max_c F(a, b, c) = n \mathbb{E} \left| \sum_{j=1}^n b_j \right| \geq \frac{n^{3/2} \gamma}{2}.$$

En particulier, pour tout choix de a , il existe b tel que $\max_c F(a, b, c) \geq \frac{n^{3/2} \gamma}{2}$.

Exercice 4.

Estimer l'intersection avec un rectangle

Let $P \subset \mathbb{Z}^2$ of size n . Our objective is to be able to answer quickly queries of the form what is the fraction of points in P that are in the rectangle $r = [a_1, b_1] \times [a_2, b_2]$? We write $r[P] = \frac{|P \cap r|}{n}$ for this fraction. We consider a simple data structure to approximate $r[P]$ efficiently for any query r . The data structure is just a random subset $S \subset P$ of size m . On query r , the estimate for $r[P]$ we output is $\frac{|S \cap r|}{m}$. The structure S defines an ε -approximation if for all queries r , we have $|r[P] - \frac{|S \cap r|}{m}| \leq \varepsilon$.

1. What m should we take to obtain an ε -approximation with probability $1 - \delta$?

☞ On va prendre un sample S dont l'espérance de la taille est m (plutôt que taille exactement m).

Pour $p \in P$, soit X_p une variable aléatoire de Bernoulli de paramètre m/n , et l'on définit S par la relation suivante : si $X_p = 1$ alors $p \in S$, et si $X_p = 0$ alors $p \notin S$. Fixons un rectangle r et soit $X(r) = \sum_{p \in R} X_p = |S \cap R|$ de telle sorte que $X(r)/m$ soit notre estimateur. Alors $\mathbb{E}[X(r)] = \sum_{p \in R} \mathbb{P}\{p \in S\} = \sum_{p \in R} \mathbb{P}\{m/n\} = mr[P]$. On peut donc appliquer Chernoff à $X(r)$ car :

$$\mathbb{P}\{|X(r)/m - r[P]| \geq \varepsilon\} = \mathbb{P}\{|X(r) - mr[P]| \geq \varepsilon m\} = \mathbb{P}\{|X(r) - \mathbb{E}[X(r)]| \geq \varepsilon/r[P] \cdot \mathbb{E}[X(r)]\} \leq 2e^{-\frac{\varepsilon^2}{2+\varepsilon} \mathbb{E}[X(r)]}.$$

avec $\varepsilon' = \varepsilon/r[P]$, ce qui donne (en utilisant $r[P] \leq 1$ pour la dernière inégalité) :

$$2e^{-\frac{\varepsilon^2}{2+\varepsilon} \mathbb{E}[X(r)]} \leq 2e^{-\frac{\varepsilon^2}{2r[P]+\varepsilon} m} \leq 2e^{-\frac{\varepsilon^2}{2+\varepsilon} m}.$$

Cette inégalité est vraie pour un rectangle r fixé, mais nous avons besoin d'une Union-Bound sur tous les rectangles. Or, il y a une infinité de rectangles possibles dans \mathbb{Z}^2 , donc nous devons être un peu plus malin. Il faut remarquer que si r et r' sont des rectangles pour lesquels $P \cap r = P \cap r'$, alors $r[P] = r'[P]$ et l'estimation sera la même, donc l'erreur sur l'un sera exactement la même que l'erreur sur l'autre. En d'autres termes, on veut trouver un certain nombre de rectangle r_1, r_2, \dots, r_k tels que pour tout rectangle r de \mathbb{Z}^2 , il existe i tel que $P \cap r = P \cap r_i$. Ainsi,

$$\mathbb{P}\{\exists r \text{ s.t. } |r[P] - X(r)| \geq \varepsilon\} \leq \sum_{i=1}^k \mathbb{P}\{|r_i[P] - X(r_i)| \geq \varepsilon\} \leq k 2e^{-\frac{\varepsilon^2}{2+\varepsilon} m}.$$

Montrons maintenant qu'on peut obtenir $k = n^4$: pour chaque 4-uplet des points de P $((x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4))$, on définit un rectangle $r_i = [x_1, x_2] \times [y_3, y_4]$. Cet ensemble de n^4 rectangles a bien la propriété demandée car si r est un rectangle quelconque, on peut

"pousser" sa limite verticale gauche le plus à droite possible jusqu'à rencontrer un point de P , auquel cas on s'arrête de "pousser". On fait de même pour les quatre côtés du rectangle (on "pousse" vers l'intérieur jusqu'à rencontrer un point de P), et on tombe sur un r_i pour lequel $r \cap P = r_i \cap P$.
 En résumé, nous voulons m tel que

$$n^4 2e^{-\frac{\varepsilon^2}{2+\varepsilon} m} \leq \delta,$$

ce qui est possible pour

$$m \geq \frac{2+\varepsilon}{\varepsilon^2} (4 \ln n + \ln 2 - \ln \delta) = \Omega(\ln n).$$

Exercice 5.

Grappe Aléatoire Bipartite

Soit $0 < p < 1$ et $n \in \mathbb{N}^*$. On définit un graphe aléatoire non orienté $H_{2n,p}$ de la manière suivante. On se donne une famille $\{X_{i,j} : 1 \leq i \leq n, n+1 \leq j \leq 2n\}$ de v.a. i.i.d. de loi de Bernoulli de paramètre p . On pose alors $H_{2n,p} = (V, E)$, avec $V = \{1, \dots, 2n\}$ et

$$E = \{(i, j) : X_{i,j} = 1\} \subset \{1, \dots, n\} \times \{n+1, \dots, 2n\}.$$

1. Quelle est la loi du nombre d'arêtes de $H_{2n,p}$?

☞ Le nombre d'arêtes de $H_{2n,p}$ suit la loi $B(n^2, p)$.

2. Quelle est l'espérance du nombre de sommets isolés de $H_{2n,p}$?

☞ Soit N le nombre de sommets isolés. Si A_i est l'événement « le sommet i est isolé », on a par linéarité de l'espérance, on a $\mathbb{E}[N] = \sum \mathbb{P}(A_i) = 2n(1-p)^n$.

3. Dans cette question on pose $p = c \log(n)/n$ pour un nombre réel $c > 0$.

1. Montrer que si $c > 1$, alors

$$\lim_{n \rightarrow \infty} \mathbb{P}(H_{2n,p} \text{ a un sommet isolé}) = 0.$$

2. Montrer que si $c < 1$, alors

$$\lim_{n \rightarrow \infty} \mathbb{P}(H_{2n,p} \text{ a un sommet isolé}) = 1.$$

☞

1. Si $c > 1$, on a $\mathbb{E}[N] = 2n \exp(n \log(1 - \frac{c \log n}{n})) \rightarrow 0$ et donc $\mathbb{P}(N \geq 1) \leq \mathbb{E}[N] \rightarrow 0$.

2. Si $c < 1$, on calcule

$$\mathbb{E}[N^2] = \sum_{i,j=1}^{2n} \mathbb{P}(A_i \cap A_j) = 2n(1-p)^n + 2n(n-1)(1-p)^{2n} + 2n^2(1-p)^{2n-1}$$

d'où il vient que $\mathbb{E}[N^2]/\mathbb{E}[N]^2$ tend vers 1. On utilise l'inégalité de Tchebychev pour conclure que

$$\mathbb{P}(N=0) = \mathbb{P}(\mathbb{E}[N] - N \geq \mathbb{E}[N]) \leq \mathbb{P}(|N - \mathbb{E}[N]| \geq \mathbb{E}[N]) \leq \frac{\text{Var}[N]}{\mathbb{E}[N]^2} = \frac{\mathbb{E}[N^2]}{\mathbb{E}[N]^2} - 1 \rightarrow 0.$$

4. Dans cette question on pose $p = 1/2$. Montrer qu'il existe une constante $C > 0$ telle que

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\text{tous les sommets de } H_{2n,p} \text{ ont un degré inférieur à } \frac{n}{2} + C\sqrt{n \log n} \right) = 1.$$

☞ Le degré d_i du sommet i suit la loi $B(n, 1/2)$. Par l'inégalité de Chernoff I, on a donc

$$\mathbb{P}(d_i \geq \frac{n}{2} + a) \leq \exp(-2a^2/n).$$

Ainsi, par la borne de l'union,

$$\mathbb{P}(\max_i d_i \geq \frac{n}{2} + a) \leq 2n \exp(-2a^2/n).$$

Cette quantité tend vers 0 si $a = C\sqrt{n \log n}$ avec $2C^2 > 1$.