



DÉTECTION DE COMMUNAUTÉS ET APPRENTISSAGE EN LIGNE

15 février 2022

ANTOINE BAK, ROMAIN FARTHOAT
sous la supervision de MATTHIEU LERASLE



REMERCIEMENTS

Merci aux organisateurs du MODAL MAP472A pour avoir rendu possible ce projet, ainsi qu'à Matthieu Lerasle pour nous avoir guidés tout au long de ce travail.

TABLE DES MATIÈRES

Remerciements	2
1 Problème du bandit manchot	4
1.1 Problème posé	4
1.2 Stratégie <i>Upper Confidence Bound</i> (UCB)	6
1.2.1 Choix de δ	6
1.3 Cas avec τ connu	8
1.3.1 Inégalité de Chernov	8
1.3.2 Calcul du regret	9
2 Apprentissage en ligne dans un graphe	13
2.1 Problème posé	13
2.2 Première solution	14
2.2.1 Schéma de preuve	14
2.2.2 Hypothèses	15
2.2.3 Évaluation de $\hat{\tau}$	15
2.2.4 Évaluation de l'erreur commise en fonction du nombre de sommets mal catégorisés	16
2.2.5 Calcul du regret	20
2.3 Amélioration possible	22
2.3.1 Conjectures sur cet algorithme	22
2.3.2 Retour sur le choix de n	23
2.3.3 Majoration de $ \mathcal{S} $	23
2.3.4 Résultats à prouver	24
2.3.5 Remarques d'ordre pratique	24

1

PROBLÈME DU BANDIT MANCHOT

1.1 PROBLÈME POSÉ

Soient $V = C_0 \sqcup C_1$ une partition de notre ensemble de sommets en deux communautés et un noeud $x \in V$ qu'on cherche à catégoriser. Ces deux communautés sont ici les « bras » du problème, selon une terminologie courante.

On notera $i^* \in \{0, 1\}$ l'indice de la communauté de notre noeud.

L'appartenance communautaire de tous les autres sommets est supposée connue dans la première partie, et seulement celle-ci.

On n'impose ici aucune contrainte sur la taille du graphe, qui peut être infini.

Notre but dans cette première partie sera de maximiser l'espérance du nombre d'arêtes découvertes en un nombre $t \in \mathbb{N}^*$ fixé de tours.

Étant donné un sommet autre que x quelconque, suppose que la probabilité qu'une arête relie x à ce sommet vaut p s'il appartient à la même communauté que x , q sinon, avec $p > q$. On veut donc d'après ce qui précède déterminer i^* , mais tout en restant efficace : on ne veut pas prendre le risque de moins gagner au final en jouant trop sur la mauvaise communauté.

On appellera μ_0 la probabilité que le sommet x soit relié à un autre sommet dans C_1 . On définit similairement μ_1 .

Conformément à ce qui précède :

$$\begin{cases} \mu_{i^*} = p \\ \mu_{1-i^*} = q \end{cases}$$

Au tour $s \in \mathbb{N}^*$ et compte tenu des informations acquises, on choisit une communauté $\pi(s)$, puis on procède à un tirage dans cette dernière, *i.e.* on y pioche un sommet — pas encore observé — pour voir s'il est relié à x ou non.

Soit, pour tout $s \in \mathbb{N}^*$, $X_{i^*,s} \sim \mathcal{B}(p)$ et $X_{1-i^*,s} \sim \mathcal{B}(q)$ avec $0 < q < p$. Ces variables représentent les tirages successifs effectués respectivement dans la communauté i^* et dans l'autre. On suppose toutes ces variables aléatoires **indépendantes**.

Posons, pour $s \in \llbracket 1, t \rrbracket$: $Y_s = X_{\pi(s), N_{\pi(s)}(s)}$ et $N_i(s)$ le nombre de fois où le bras $i \in \{0, 1\}$ a été choisi au cours des s premiers instants.

On cherche une stratégie $\pi : \llbracket 1, t \rrbracket \longrightarrow \{0, 1\}$ optimale, *i.e.* minimisant la fonction de regret suivante :

$$R_t = tp - \mathbb{E} \left[\sum_{s=1}^t Y_s \right]$$

$\pi(s)$ peut bien sûr dépendre des résultats passés et même être randomisé, *i.e.* défini comme suit :

$$\pi(s) = \psi_s(U_1, Y_1, \dots, U_{s-1}, Y_{s-1}, U_s)$$

où les U_i sont i.i.d. de loi $\mathcal{U}([0, 1])$ et indépendantes des $X_{i,s}$.

Remarque. Il n'y a pas de raison pour que $\pi(s)$ ne dépende pas de t , nous avons seulement allégé les notations.

On notera enfin à travers l'ensemble du présent rapport :

$$\Delta = p - q \qquad \tau = \frac{p + q}{2}$$

Revenons au risque. La formule de Wald permet de montrer qu'on a :

$$R_t = \Delta \mathbb{E} [N_{1-i^*}(t)]$$

et que de plus, en définissant :

$$\forall s \in \llbracket 1, t \rrbracket, \forall i \in \{0, 1\}, \quad \hat{\mu}_i(s) = \frac{1}{N_i(s)} \sum_{k=1}^{N_i(s)} X_{i,k}$$

on trouve des estimateurs sans biais des μ_i :

$$\forall s \in \llbracket 1, t \rrbracket, \forall i \in \{0, 1\}, \quad \mathbb{E}[\hat{\mu}_i(s)] = \mu_i$$

1.2 STRATÉGIE *UPPER CONFIDENCE BOUND* (UCB)

Soit δ une fonction décroissante dont on déterminera les valeurs plus tard. On pose :

$$\Omega = \bigcap_{\substack{1 \leq s \leq t \\ i \in \{0,1\}}} \left(\mu_i \geq \hat{\mu}_i(s) - \delta(s) \cap \mu_{i^*} \leq \hat{\mu}_{i^*}(s) + \delta(s) \right)$$

Remarque. Il n'y a pas non plus de raison pour que δ ne dépende pas de t .

L'évènement Ω correspond au cas où « tout va bien », *i.e.* où, peu importe le rang s auquel on s'intéresse, aucun des deux estimateurs n'est « trop grand », et où $\hat{\mu}_{i^*}(s)$ n'est pas « trop petit ». $\delta(s)$ représente ici l'écart à l'espérance que l'on s'autorise avant de taxer un estimateur de « trop grand » ou « trop petit ».

On peut décider de choisir i^* constamment à compter du rang $s_0 \in \mathbb{N}^*$ si pour $s_0 \leq s \leq t$:

$$\hat{\mu}_{1-i^*}(s) + \delta(s) < \hat{\mu}_{i^*}(s) - \delta(s) \quad (1)$$

Or quitte à supposer Ω satisfait, on a pour $1 \leq s \leq t$:

$$\begin{cases} \mu_{i^*} - \delta(s) \leq \hat{\mu}_{i^*}(s) \\ \hat{\mu}_{1-i^*}(s) \leq \mu_{1-i^*} + \delta(s) \end{cases}$$

d'où la condition suffisante (par décroissance de δ) pour (1) :

$$\mu_{1-i^*} + 2\delta(s_0) \leq \mu_{i^*}$$

Cette condition est satisfaite si $s_0 \geq \delta^{-1}\left(\frac{\Delta}{2}\right)$, δ^{-1} étant la fonction inverse généralisée de δ .

Cela entraîne $\mathbb{E}[N_{1-i^*}(t)|\Omega] \leq \delta^{-1}\left(\frac{\Delta}{2}\right)$ et :

$$R_t \leq \Delta \delta^{-1}\left(\frac{\Delta}{2}\right) + t \Delta \mathbb{P}(\Omega^C)$$

Restent à déterminer $\delta^{-1}\left(\frac{\Delta}{2}\right)$ et $\mathbb{P}(\Omega^C)$.

1.2.1 • CHOIX DE δ

Comme :

$$\Omega \supset \bigcap_{1 \leq s \leq t} \left(|\hat{\mu}_0(s) - \mu_0| > \delta(s) \cap |\hat{\mu}_1(s) - \mu_1| > \delta(s) \right)$$

on va utiliser l'inégalité de Hoeffding.

Pour $s \in \llbracket 1, t \rrbracket$ et $N_0(s)$ **donnés** :

$$\begin{aligned}
& \mathbb{P}\left(|\hat{\mu}_0(s) - \mu_0| > \delta(s) \cap |\hat{\mu}_1(s) - \mu_1| > \delta(s)\right) \\
& \leq \mathbb{P}\left(N_0(s)|\hat{\mu}_0(s) - \mu_0| > N_0(s)\delta(s) \cap N_1(s)|\hat{\mu}_1(s) - \mu_1| > N_1(s)\delta(s)\right) \\
& \leq \mathbb{P}\left(\left|\sum_{k=1}^{N_0(s)} X_{0,k} - N_0(s)\mu_0\right| > N_0(s)\delta(s) \cap \left|\sum_{k=1}^{N_1(s)} X_{1,k} - N_1(s)\mu_1\right| > N_1(s)\delta(s)\right) \\
& \leq \mathbb{P}\left(\left|\sum_{k=1}^{N_0(s)} X_{0,k} + \sum_{k=1}^{N_1(s)} X_{1,k} - (N_0(s)\mu_0 + N_1(s)\mu_1)\right| > (N_0(s) + N_1(s))\delta(s)\right) \\
& \leq 2 \exp\left(-2(N_0(s) + N_1(s))\delta(s)^2\right) \\
& \leq 2 \exp\left(-2s\delta(s)^2\right)
\end{aligned}$$

On en déduit, indépendamment de la valeur de $N_0(s)$:

$$\mathbb{P}\left(|\hat{\mu}_0(s) - \mu_0| > \delta(s) \cap |\hat{\mu}_1(s) - \mu_1| > \delta(s)\right) \leq 2 \exp\left(-2s\delta(s)^2\right)$$

Soit $\eta > 0$. Avec $\delta(s) = \sqrt{\frac{\log(2t/\eta)}{2s}}$ on obtient $\mathbb{P}\left(|\hat{\mu}_i(s) - \mu_i| > \delta(s)\right) \leq \frac{\eta}{t}$ et :

$$\begin{aligned}
\mathbb{P}\left(\Omega^C\right) & \leq \mathbb{P}\left(\bigcup_{1 \leq s \leq t} \left(|\hat{\mu}_0(s) - \mu_0| > \delta(s) \cap |\hat{\mu}_1(s) - \mu_1| > \delta(s)\right)\right) \\
& \leq \sum_{1 \leq s \leq t} \mathbb{P}\left(|\hat{\mu}_0(s) - \mu_0| > \delta(s) \cap |\hat{\mu}_1(s) - \mu_1| > \delta(s)\right) \\
& \leq \eta
\end{aligned}$$

Enfin, $\eta = \frac{1}{t}$ donne la majoration :

$$R_t \leq \Delta + \frac{2 \log(2t^2)}{\Delta}$$

1.3 CAS AVEC τ CONNU

Une idée naturelle consiste à observer la moyenne empirique des tirages pour chaque communauté et à la comparer avec la valeur de τ tant qu'aucune communauté n'a été épuisée :

- si les deux sont inférieures à τ , ou qu'aucun sommet n'a encore été tiré, on effectue un tirage supplémentaire dans chaque communauté ;
- sinon on tire un sommet dans la communauté ayant la moyenne empirique la plus grande.

On suppose $\Delta \leq 2q$ pour pouvoir plus tard appliquer l'inégalité de Chernov.

1.3.1 • INÉGALITÉ DE CHERNOV

Théorème 1. Soit $(A_i)_{i \geq 1}$ une suite de variables aléatoires i.i.d. de loi $\mathcal{B}(p)$ et $t \in \mathbb{N}^*$.

Soit $S_t = \sum_{i=1}^t A_i$, de sorte que $\mathbb{E}[S_t] = tp$. On note également $\bar{A}_t = \frac{1}{t} \sum_{i=1}^t A_i$.

Si $\delta \in [0, 1]$, alors :

$$\begin{cases} \mathbb{P}(S_t \geq tp(1 + \delta)) \leq \exp\left(-\frac{tp\delta^2}{3}\right) \\ \mathbb{P}(S_t \leq tp(1 - \delta)) \leq \exp\left(-\frac{tp\delta^2}{3}\right) \end{cases}$$

Si $\delta \geq 1$, on a quand même :

$$\mathbb{P}(S_t \geq tp(1 + \delta)) \leq \exp\left(-\frac{tp\delta}{3}\right)$$

Remarque. Cette inégalité a pour intérêt d'obtenir de meilleures majorations que l'inégalité de Hoeffding dans les cas où p est faible :

Si $\varepsilon \leq p$:

$$\begin{aligned} \mathbb{P}(|\bar{A}_t - p| \geq \varepsilon) &\leq \mathbb{P}\left(S_t \geq pl \left(1 + \frac{\varepsilon}{p}\right)\right) + \mathbb{P}\left(S_t \leq pl \left(1 - \frac{\varepsilon}{p}\right)\right) \\ &\leq 2 \exp\left(-\frac{l\varepsilon^2}{3p}\right) \end{aligned}$$

Démonstration. On applique l'inégalité de Markov à $e^{\lambda S_t}$ pour $\lambda > 0$:

$$\begin{aligned} \mathbb{P}(S_t \geq tp(1 + \delta)) &= \mathbb{P}(e^{\lambda S_t} \geq e^{\lambda tp(1 + \delta)}) \\ &\leq \mathbb{E}[e^{\lambda S_t}] e^{-\lambda tp(1 + \delta)} \\ &= \exp(t \log(1 + p(e^\lambda - 1)) - \lambda tp(1 + \delta)) \end{aligned}$$

D'une part $\lambda = \log(1 + \delta)$ donne :

$$\begin{aligned}\mathbb{P}(S_t \geq tp(1 + \delta)) &\leq \exp(t \log(1 + p\delta) - tp(1 + \delta) \log(1 + \delta)) \\ &\leq \exp(-tp((1 + \delta) \log(1 + \delta) - \delta))\end{aligned}$$

de même :

$$\begin{aligned}\mathbb{P}(S_t \leq tp(1 - \delta)) &\leq \mathbb{E}[e^{-\lambda S_t}] e^{\lambda tp(1 - \delta)} \\ &= \exp(t \log(1 + p(e^{-t} - 1)) + \lambda tp(1 - \delta))\end{aligned}$$

D'autre part $\lambda = -\log(1 - \delta)$ donne :

$$\mathbb{P}(S_t \leq tp(1 - \delta)) \leq \exp(-tp((1 - \delta) \log(1 - \delta) + \delta))$$

Avec une rapide étude de fonction, on trouve que pour $\delta \in [-1, 1]$:

$$(1 + \delta) \log(1 + \delta) - \delta \geq \frac{\delta^2}{3}$$

ce qui achève la démonstration. □

1.3.2 • CALCUL DU REGRET

Intéressons-nous à présent à la probabilité de se tromper au s^e tour avec cette stratégie :

$$\mathbb{P}(\pi(s) = 1 - i^*) = \mathbb{P}(\pi(s) = 1 - i^* \cap \hat{\mu}_{1-i^*}(s) > \tau) + \mathbb{P}(\pi(s) = 1 - i^* \cap \hat{\mu}_{1-i^*}(s) \leq \tau)$$

Or :

$$\begin{aligned}
& \sum_{s=3}^t \mathbb{P}(\pi(s) = 1 - i^* \cap \hat{\mu}_{1-i^*}(s) > \tau) \\
&= \sum_{s=3}^t \mathbb{P}\left(N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{1-i^*}(s) - \mu_{1-i^*} > \frac{\Delta}{2}\right) \\
&= \mathbb{E}\left[\sum_{s=3}^t \mathbf{1}_{N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{1-i^*}(s) - \mu_{1-i^*} > \frac{\Delta}{2}}\right] \\
&= \mathbb{E}\left[\sum_{\substack{3 \leq s \leq t \\ N_{1-i^*}(s) > N_{1-i^*}(s-1)}} \mathbf{1}_{\hat{\mu}_{1-i^*}(s) - \mu_{1-i^*} > \frac{\Delta}{2}}\right] \\
&= \mathbb{E}\left[\sum_{l=1}^{N_{1-i^*}(t)} \mathbf{1}_{\frac{1}{l} \sum_{k=1}^l X_{1-i^*,k} - \mu_{1-i^*} > \frac{\Delta}{2}}\right] \\
&\leq \mathbb{E}\left[\sum_{l=1}^t \mathbf{1}_{\sum_{k=1}^l X_{1-i^*,k} - l\mu_{1-i^*} > l\frac{\Delta}{2}}\right] \\
&= \sum_{l=1}^t \mathbb{P}\left(\sum_{k=1}^l X_{1-i^*,k} > lq \left(1 + \frac{\Delta}{2q}\right)\right) \\
&\leq \sum_{l=1}^t \exp\left(-\frac{l\Delta^2}{12q}\right) \quad \text{par Chernov} \\
&\leq \frac{1}{1 - \exp\left(-\frac{\Delta^2}{12q}\right)} \\
&\leq \frac{12q}{\Delta^2}
\end{aligned}$$

De même, si :

$$\begin{cases} \pi(s) = 1 - i^* \\ \hat{\mu}_{1-i^*}(s) \leq \tau \end{cases}$$

on a nécessairement $\hat{\mu}_{i^*}(s) \leq \hat{\mu}_{1-i^*}(s) \leq \tau$ car c'est $1 - i^*$ qui est choisi au s^e tour.

Il vient alors similairement que :

$$\begin{aligned}
& \sum_{s=3}^t \mathbb{P}(\pi(s) = 1 - i^* \cap \hat{\mu}_{1-i^*}(s) \leq \tau) \\
& \leq \sum_{s=3}^t \mathbb{P}(\pi(s) = 1 - i^* \cap \hat{\mu}_{i^*}(s) \leq \tau) \\
& \leq \sum_{s=3}^t \mathbb{P}\left(N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{i^*}(s) - \mu_{i^*} \leq -\frac{\Delta}{2}\right) \\
& = \mathbb{E}\left[\sum_{s=3}^t \mathbf{1}_{N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{i^*}(s) - \mu_{i^*} < -\frac{\Delta}{2}}\right] \\
& = \mathbb{E}\left[\sum_{\substack{3 \leq s \leq t \\ N_{1-i^*}(s) > N_{1-i^*}(s-1)}} \mathbf{1}_{\hat{\mu}_{i^*}(s) - \mu_{i^*} < -\frac{\Delta}{2}}\right] \\
& = \mathbb{E}\left[\sum_{l=1}^{N_{1-i^*}(t)} \mathbf{1}_{\frac{1}{l} \sum_{k=1}^l X_{i^*,k} - \mu_{i^*} < -\frac{\Delta}{2}}\right] \\
& \leq \mathbb{E}\left[\sum_{l=1}^t \mathbf{1}_{\sum_{k=1}^l X_{i^*,k} - l\mu_{i^*} < -l\frac{\Delta}{2}}\right] \\
& = \sum_{l=1}^t \mathbb{P}\left(\sum_{k=1}^l X_{i^*,k} < lp\left(1 - \frac{\Delta}{2p}\right)\right) \\
& \leq \sum_{l=1}^t \exp\left(-\frac{l\Delta^2}{12p}\right) \quad \text{par Chernov} \\
& \leq \frac{1}{1 - \exp\left(-\frac{\Delta^2}{12p}\right)} \\
& \leq \frac{12p}{\Delta^2}
\end{aligned}$$

Par ailleurs comme $N_{1-i^*}(t) = \sum_{s=1}^t \mathbf{1}_{\pi(s)=1-i^*}$, on trouve :

$$\begin{aligned}
 \mathbb{E}[N_{1-i^*}(t)] &= \mathbb{E}\left[\sum_{s=1}^t \mathbf{1}_{\pi(s)=1-i^*}\right] \\
 &= \sum_{s=1}^t \mathbb{E}\left[\mathbf{1}_{\pi(s)=1-i^*}\right] \\
 &= \sum_{s=1}^t \mathbb{P}(\pi(s) = 1 - i^*) \\
 &\leq \frac{12(p+q)}{\Delta^2} \\
 &= 24 \frac{\tau}{\Delta^2}
 \end{aligned}$$

et finalement :

$$R_t \leq \Delta + \frac{24\tau}{\Delta}$$

Le premier terme de la somme vient des deux tirages d'initialisation de l'algorithme.

Remarque. En utilisant l'inégalité de Hoeffding plutôt que celle de Chernov on peut obtenir :

$$R_t \leq \Delta + \frac{4}{\Delta}$$

Cette dernière inégalité a le mérite d'être valable même si $\Delta > 2q$.

2

APPRENTISSAGE EN LIGNE DANS UN GRAPHE

2.1 PROBLÈME POSÉ

On reprend un formalisme similaire : soit $V = C_0 \sqcup C_1$ une partition des sommets d'un graphe en deux communautés telles que si i et j sont des sommets quelconque de la même communauté ils ont une probabilité p d'être reliés, contre q s'ils sont de communautés différentes.

On définit naturellement les variables aléatoires indépendantes :

$$X_{ij} \sim \begin{cases} \mathcal{B}(p) & \text{si } i \text{ et } j \text{ appartiennent à la même communauté} \\ \mathcal{B}(q) & \text{sinon} \end{cases}$$

Le problème posé est proche de celui du bandit manchot, où les « machines » sont les arêtes du graphe, et où il n'est pas possible de tirer deux fois le même levier.

Pour un nombre donné de tirages $T \in \mathbb{N}^*$, on veut jouer de manière à maximiser le nombre d'arêtes découvertes. Cette fois-ci en revanche, **on a a priori aucune connaissance sur l'appartenance communautaire de nos sommets.**

Le seul paramètre d'entrée de l'algorithme est l'entier $T \in \mathbb{N}^*$ qui désigne le nombre de tirages à effectuer au total.

2.2 PREMIÈRE SOLUTION

1. On tire un ensemble \mathcal{N} de n sommets de manière équilibrée, numérotés de 1 à n , puis :
 - on tire les arêtes correspondant à chaque paire du noyau ;
 - on évalue $\hat{\tau} = \binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} X_{ij}$;
 - on utilise un algorithme de détection de communautés dans le sous-graphe \mathcal{N} .
2. On se propose de classer un ensemble \mathcal{S} de sommets de $G \setminus \mathcal{N}$.
 Tant que nous n'avons pas effectué T tirages, on pioche $x \in V \setminus \mathcal{N}$, et on réutilise la stratégie avec τ connu en l'approximant par $\hat{\tau}$, le tout sur les communautés \hat{C}_0 et \hat{C}_1 :
 - On tire un sommet dans chaque \hat{C} .
 - Tant que ni \hat{C}_0 ni \hat{C}_1 n'a été intégralement visitée :
 - Si les $\hat{\mu}$ sont tous deux inférieurs à $\hat{\tau}$, on tire un sommet supplémentaire dans chaque \hat{C} .
 - Sinon, on tire un sommet de la communauté associée au $\hat{\mu}$ le plus élevé.

2.2.1 • SCHÉMA DE PREUVE

Les notations utilisées ici seront introduites rigoureusement lors des développements suivants. On ne vise ici qu'à donner une idée globale de la direction des démonstrations qui suivront.

On montre tout d'abord que pourvu que notre noyau ait une composition équilibrée, *i.e.* qu'aucune communauté réelle n'y soit surreprésentée :

$$\left| n_0 - \frac{n}{2} \right| \leq n\varepsilon_n$$

on peut utiliser la méthode de classification du cas τ connu en l'estimant de manière satisfaisante, soit encore :

$$|\hat{\tau} - \tau| \leq \Delta\varepsilon_\tau$$

On montre ensuite que pour chaque sommet donné, notre procédé de classification aboutit à un nombre borné d'erreurs, pourvu que la catégorisation du noyau soit suffisamment bonne, c'est-à-dire :

$$|\hat{C}_0 \cap C_1| + |\hat{C}_1 \cap C_0| \leq n\varepsilon_C$$

Enfin, le calcul du regret permet d'évaluer cet algorithme. On montre que la grande probabilité de la conjonction de ces trois conditions (notée Ω) rend possible un regret asymptotique en $O(T^{2/3})$ au lieu de $O(T)$ en tirant simplement des arêtes au hasard.

2.2.2 • HYPOTHÈSES

On suppose que $\Delta \geq \frac{C}{n\varepsilon_C^2}$

Pour pouvoir utiliser l'inégalité de Chernov, on suppose que $\sqrt{35}q \geq \Delta$.

2.2.3 • ÉVALUATION DE $\hat{\tau}$

Posons $n_0 = |\mathcal{N} \cap C_0|$ et $n_1 = |\mathcal{N} \cap C_1| = n - n_0$.

Remarque. Noter que n_0 et n_1 dépendent du noyau, que l'on sélectionne aléatoirement. Ce sont donc **des variables aléatoires**.

On a quitte à renuméroter les sommets :

$$\binom{n}{2}\hat{\tau} = \sum_{1 \leq i < j \leq n_0} X_{ij} + \sum_{n_0+1 \leq i < j \leq n} X_{ij} + \sum_{1 \leq i \leq n_0, n_0+1 \leq j \leq n} X_{ij}$$

où les X_{ij} des deux premières sommes suivent une loi $\mathcal{B}(p)$, alors que celles de la dernière somme suivent une loi $\mathcal{B}(q)$.

On peut donc écrire :

$$\binom{n}{2}\hat{\tau} \sim B_1 + B_2$$

où

$$B_1 \sim \mathcal{B}\left(\binom{n_0}{2} + \binom{n_1}{2}, p\right)$$

et

$$B_2 \sim \mathcal{B}(n_0 n_1, q)$$

sont deux variables aléatoires indépendantes pour toute valeur de n_0 **donnée**.

On en déduit l'expression de l'espérance de $\hat{\tau}$:

$$\binom{n}{2}\mathbb{E}[\hat{\tau}|n_0] = p\left(\binom{n_0}{2} + \binom{n_1}{2}\right) + qn_0 n_1$$

Soit $n\varepsilon_n = \sqrt{\frac{\Delta}{32\tau}\binom{n}{2}}$

Si $\left|n_0 - \frac{n}{2}\right| \leq n\varepsilon_n$, on a :

$$\begin{cases} \left|n_0 n_1 - \frac{1}{2} \binom{n}{2}\right| \leq \frac{\Delta}{32\tau} \binom{n}{2} \\ \left|\binom{n_0}{2} + \binom{n_1}{2} - \frac{1}{2} \binom{n}{2}\right| \leq \frac{\Delta}{32\tau} \binom{n}{2} \end{cases}$$

En utilisant l'inégalité de Serfling [1] :

$$\mathbb{P}\left(\left|n_0 - \frac{n}{2}\right| > n\varepsilon_n\right) \leq 2 \exp(-2n\varepsilon_n^2)$$

On aimerait à présent avoir une estimation suffisamment précise de τ pour qu'utiliser les résultats de la première partie soit pertinent.

Plus concrètement, on désire avoir $|\hat{\tau} - \tau| \leq \Delta\varepsilon_\tau$.

On peut appliquer l'inégalité de Hoeffding pour trouver :

$$\mathbb{P}\left(|\hat{\tau} - \tau| > \Delta\varepsilon_\tau \mid \left|n_0 - \frac{n}{2}\right| \leq n\varepsilon_n\right) \leq 2 \exp\left(-2\binom{n}{2}\Delta^2(\varepsilon_\tau - \varepsilon_n)^2\right)$$

2.2.4 • ÉVALUATION DE L'ERREUR COMMISE EN FONCTION DU NOMBRE DE SOMMETS MAL CATÉGORISÉS

On peut montrer que s'il y a m sommets mal catégorisés dans \hat{C}_i alors :

$$\mathbb{E}[\hat{\mu}_i(s) | \hat{C}_i] = \frac{\mu_i(|\hat{C}_i| - m) + \mu_{1-i}m}{|\hat{C}_i|}$$

Or avec une probabilité supérieure à $1 - e^{-N}$ on trouve des catégories estimées telles que :

$$|\hat{C}_0 \cap C_1| + |\hat{C}_1 \cap C_0| \leq n\varepsilon_C$$

On définit alors :

$$\Omega = \left(\left|n_0 - \frac{n}{2}\right| \leq n\varepsilon_n\right) \cap (|\hat{\tau} - \tau| \leq \Delta\varepsilon_\tau) \cap (|\hat{C}_0 \cap C_1| + |\hat{C}_1 \cap C_0| \leq n\varepsilon_C)$$

On peut alors utiliser une démonstration analogue à celle du cas avec τ connu.

On reprend là encore des notations similaires. À supposer que l'on s'occupe d'un sommet $x \in V \setminus \mathcal{N}$ appartenant à C_{i^*} , $\pi(s)$ désigne l'indice de la communauté choisie lors du s^e tirage effectué pour classer x . Au cours des s premiers tours de cette classification, la communauté i a été choisie $N_i(s)$ fois.

Par la formule des probabilités totales :

$$\begin{aligned}
 \mathbb{P}(\pi(s) = 1 - i^* | \Omega) &= \mathbb{P}(\pi(s) = 1 - i^* \cap \hat{\mu}_{1-i^*}(s) \leq \hat{\tau} | \Omega) \\
 &\quad + \mathbb{P}(\pi(s) = 1 - i^* \cap \hat{\mu}_{1-i^*}(s) > \hat{\tau} | \Omega) \\
 &= \mathbb{P}(N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{1-i^*}(s) > \hat{\tau} | \Omega) \\
 &\quad + \mathbb{P}(N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{1-i^*}(s) \leq \hat{\tau} | \Omega)
 \end{aligned}$$

Puis comme dans la démonstration du cas τ connu, on note $(Y_{i,k})_{1 \leq k \leq |\hat{C}_i|}$ les variables aléatoires représentant les tirages d'arêtes entre x et les sommets de \hat{C}_i pris dans un ordre aléatoire et $\bar{Y}_{i,l} = \frac{1}{l} \sum_{k=1}^l Y_{i,k}$:

$$\begin{aligned}
 \sum_{s=3}^t \mathbb{P}(N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{1-i^*}(s) > \hat{\tau} | \Omega) \\
 \leq \sum_{l=1}^t \mathbb{P}\left(\bar{Y}_{1-i^*,l} - \mathbb{E}[\bar{Y}_{1-i^*,l}] > \hat{\tau} - \tau + \tau - \mathbb{E}[\bar{Y}_{1-i^*,l}] \mid \Omega\right) \\
 \leq \sum_{l=1}^t \mathbb{P}\left(\bar{Y}_{1-i^*,l} - \mathbb{E}[\bar{Y}_{1-i^*,l}] > \tau - \mathbb{E}[\bar{Y}_{1-i^*,l}] - |\hat{\tau} - \tau| \mid \Omega\right)
 \end{aligned}$$

On se convainc rapidement que le nombre de sommets mal catégorisés sur les l premiers choisis dans \hat{C}_{1-i^*} suit une loi hypergéométrique :

$$\mathcal{H}\left(l, \frac{|\hat{C}_{1-i^*} \cap C_{i^*}|}{|\hat{C}_{1-i^*}|}, |\hat{C}_{1-i^*}|\right)$$

Rappelons que sous Ω , $|\hat{C}_{1-i^*} \cap C_{i^*}| \leq n\varepsilon_C$ et :

$$\begin{aligned}
 |\hat{C}_{1-i^*}| &= n_{1-i^*} + |\hat{C}_{1-i^*} \cap C_{i^*}| - |\hat{C}_{i^*} \cap C_{1-i^*}| \\
 &\geq n_{1-i^*} - |\hat{C}_{i^*} \cap C_{1-i^*}| \\
 &\geq \frac{n}{2} - n\varepsilon_n - n\varepsilon_C \\
 &= n \frac{1 - 2\varepsilon_n - 2\varepsilon_C}{2}
 \end{aligned}$$

La part $\frac{|\hat{C}_{1-i^*} \cap C_{i^*}|}{|\hat{C}_{1-i^*}|}$ de sommets mal classés dans \hat{C}_{1-i^*} valant au plus

$$\frac{n\varepsilon_C}{n \frac{1 - 2\varepsilon_n - 2\varepsilon_C}{2}} = \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C}$$

sous Ω , on a :

$$\mathbb{E} \left[\bar{Y}_{1-i^*,l} \middle| \Omega \right] \leq \mu_{i^*} \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} + \mu_{1-i^*} \left(1 - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} \right)$$

D'où :

$$\tau - \mathbb{E} \left[\bar{Y}_{1-i^*,l} \middle| \Omega \right] \geq \Delta \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} \right)$$

Et, en notant t le nombre de tours nécessaire à la classification de x :

$$\begin{aligned} \sum_{s=3}^t \mathbb{P} \left(N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{1-i^*}(s) > \hat{\tau} \middle| \Omega \right) \\ \leq \sum_{l=1}^t \mathbb{P} \left(\bar{Y}_{1-i^*,l} - \mathbb{E} \left[\bar{Y}_{1-i^*,l} \middle| \Omega \right] > \tau - \mathbb{E} \left[\bar{Y}_{1-i^*,l} \middle| \Omega \right] - |\hat{\tau} - \tau| \middle| \Omega \right) \\ \leq \sum_{l=1}^t \mathbb{P} \left(\bar{Y}_{1-i^*,l} - \mathbb{E} \left[\bar{Y}_{1-i^*,l} \middle| \Omega \right] > \Delta \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right) \middle| \Omega \right) \end{aligned}$$

Par l'inégalité de Chernov pour $1 \leq l \leq t$:

$$\begin{aligned} \mathbb{P} \left(\bar{Y}_{1-i^*,l} - \mathbb{E} \left[\bar{Y}_{1-i^*,l} \middle| \Omega \right] > \Delta \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right) \middle| \Omega \right) \\ \leq \exp \left(- \frac{l \Delta^2 \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right)^2}{3 \mathbb{E} \left[\bar{Y}_{1-i^*,l} \middle| \Omega \right]} \right) \end{aligned}$$

On obtient une majoration de la somme comme précédemment :

$$\begin{aligned} \sum_{l=1}^t \mathbb{P} \left(\bar{Y}_{1-i^*,l} - \mathbb{E} \left[\bar{Y}_{1-i^*,l} \middle| \Omega \right] > \Delta \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right) \middle| \Omega \right) \\ \leq \frac{3p}{\Delta^2 \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right)^2} \end{aligned}$$

Similairement :

$$\frac{|\hat{C}_{i^*} \cap C_{1-i^*}|}{|\hat{C}_{i^*}|} \leq \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C}$$

Donc comme $\mu_{i^*} > \mu_{1-i^*}$:

$$\mathbb{E} \left[\bar{Y}_{i^*,l} \middle| \Omega \right] \geq \mu_{i^*} \left(1 - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} \right) + \mu_{1-i^*} \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C}$$

On a :

$$\begin{aligned}
& \sum_{s=3}^t \mathbb{P} \left(N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{i^*}(s) \leq \hat{\tau} \middle| \Omega \right) \\
& \leq \sum_{l=1}^t \mathbb{P} \left(\bar{Y}_{i^*,l} - \mathbb{E} [\bar{Y}_{i^*,l} | \Omega] \leq \hat{\tau} - \tau + \tau - \mathbb{E} [\bar{Y}_{i^*,l} | \Omega] \middle| \Omega \right) \\
& \leq \sum_{l=1}^t \mathbb{P} \left(\bar{Y}_{i^*,l} - \mathbb{E} [\bar{Y}_{i^*,l} | \Omega] \leq \tau - \mathbb{E} [\bar{Y}_{i^*,l} | \Omega] + |\hat{\tau} - \tau| \middle| \Omega \right) \\
& \leq \sum_{l=1}^t \mathbb{P} \left(\bar{Y}_{1-i^*,l} - \mathbb{E} [\bar{Y}_{1-i^*,l} | \Omega] \leq -\Delta \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right) \middle| \Omega \right) \\
& \leq \frac{3p}{\Delta^2 \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right)^2}
\end{aligned}$$

Et enfin le résultat désiré :

$$\begin{aligned}
& \sum_{s=3}^t \mathbb{P} \left(\pi(s) = 1 - i^* \middle| \Omega \right) \\
& = \sum_{s=3}^t \left(\mathbb{P} \left(\pi(s) = 1 - i^*, \hat{\mu}_{1-i^*}(s) > \hat{\tau} \middle| \Omega \right) \right. \\
& \quad \left. + \mathbb{P} \left(\pi(s) = 1 - i^*, \hat{\mu}_{1-i^*}(s) \leq \hat{\tau} \middle| \Omega \right) \right) \\
& \leq \sum_{s=3}^t \left(\mathbb{P} \left(N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{1-i^*}(s) > \hat{\tau} \middle| \Omega \right) \right. \\
& \quad \left. + \mathbb{P} \left(N_{1-i^*}(s) > N_{1-i^*}(s-1) \cap \hat{\mu}_{i^*}(s) \leq \hat{\tau} \middle| \Omega \right) \right) \\
& \leq \sum_{l=1}^t \left(\mathbb{P} \left(\bar{Y}_{1-i^*,l} > \hat{\tau} \middle| \Omega \right) + \mathbb{P} \left(\bar{Y}_{i^*,l} \leq \hat{\tau} \middle| \Omega \right) \right) \\
& \leq \sum_{l=1}^t \left(\mathbb{P} \left(\bar{Y}_{1-i^*,l} - \mathbb{E} [\bar{Y}_{i^*,l} | \Omega] > \Delta \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right) \middle| \Omega \right) \right. \\
& \quad \left. + \mathbb{P} \left(\bar{Y}_{i^*,l} - \mathbb{E} [\bar{Y}_{i^*,l} | \Omega] \leq -\Delta \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right) \middle| \Omega \right) \right) \\
& \leq \frac{6p}{\Delta^2 \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right)^2}
\end{aligned}$$

2.2.5 • CALCUL DU REGRET

Sur la première étape de l'algorithme avec la construction du graphe, on commet au plus $\binom{n}{2}$ erreurs, d'où un premier terme :

$$R_{T,1} = \Delta \binom{n}{2}$$

Par ailleurs le regret de la 2^e étape $R_{T,2}$ vérifie :

$$\mathbb{E}[R_{T,2}] = \mathbb{E}[R_{T,2}|\Omega]\mathbb{P}(\Omega) + \mathbb{E}[R_{T,2}|\Omega^C]\mathbb{P}(\Omega^C) \leq \mathbb{E}[R_{T,2}|\Omega] + \Delta T \mathbb{P}(\Omega^C)$$

Pour atteindre les T tirages, comme chaque sommet x choisi dans le graphe sera affecté à au moins $\min(|\hat{C}_0|, |\hat{C}_1|)$ sommets lors de la seconde étape.

On classe donc, en plus du noyau, un nombre de sommets vérifiant :

$$|\mathcal{S}| \leq \frac{T}{\min(|\hat{C}_0|, |\hat{C}_1|)}$$

Or sous Ω , on a :

$$\min(|\hat{C}_0|, |\hat{C}_1|) \geq \frac{n}{2} - n(\varepsilon_n + \varepsilon_C)$$

On peut donc majorer le nombre de sommets à classer en plus du noyau.

De plus, on a, pour chacun de ces sommets, une espérance du nombre d'erreurs majorée par :

$$1 + \frac{6p}{\Delta^2 \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right)^2}$$

Cela donne :

$$\mathbb{E}[R_{T,2}|\Omega] \leq \left(\Delta + \frac{6p}{\Delta \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right)^2} \right) \frac{T}{n \left(\frac{1}{2} - (\varepsilon_n + \varepsilon_C) \right)}$$

Par ailleurs des résultats déjà établis donnent :

$$\mathbb{P}(\Omega^C) \leq 2 \exp \left(-2 \binom{n}{2} \Delta^2 (\varepsilon_\tau - \varepsilon_n)^2 \right) + e^{-n} + 2 \exp \left(-2n\varepsilon_n^2 \right)$$

On obtient enfin la borne suivante sur le regret global :

$$\begin{aligned} R_T &\leq \Delta \binom{n}{2} + \left(\Delta + \frac{6p}{\Delta \left(\frac{1}{2} - \frac{2\varepsilon_C}{1 - 2\varepsilon_n - 2\varepsilon_C} - \varepsilon_\tau \right)^2} \right) \frac{T}{n \left(\frac{1}{2} - (\varepsilon_n + \varepsilon_C) \right)} \\ &\quad + \Delta T \left(2 \exp \left(-2 \binom{n}{2} \Delta^2 (\varepsilon_\tau - \varepsilon_n)^2 \right) + e^{-n} + 2 \exp \left(-2n\varepsilon_n^2 \right) \right) \end{aligned}$$

On a posé :

$$n\varepsilon_n = \sqrt{\frac{\Delta}{32\tau} \binom{n}{2}}$$

En prenant par exemple : $\varepsilon_\tau = \frac{1}{4}$ et $\varepsilon_C = \frac{1}{16}$, on trouve :

$$R_T \leq \Delta \binom{n}{2} + \left(\Delta + \frac{210p}{\Delta} \right) \frac{4T}{n} + \Delta T \left(2 \exp \left(- \binom{n}{2} \frac{\Delta^2}{32} \right) + e^{-n} + 2 \exp \left(- \frac{\Delta(n-1)}{32\tau} \right) \right)$$

On choisit alors :

$$n = \frac{p^{1/3}}{\Delta^{2/3}} T^{1/3}$$

Ce qui donne :

$$R_T = O \left(\frac{p^{2/3}}{\Delta^{1/3}} T^{2/3} \right)$$

2.3 AMÉLIORATION POSSIBLE

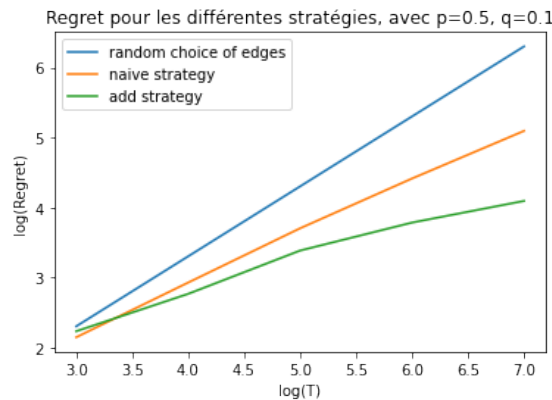
1. On tire un ensemble \mathcal{N} de n sommets de manière équilibrée, numérotés de 1 à n , puis :
 - on tire les arêtes correspondant à chaque paire du noyau ;
 - on évalue $\hat{\tau} = \binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} X_{ij}$;
 - on utilise un algorithme de détection de communautés dans le sous-graphe \mathcal{N} .
2. On se propose de classer un ensemble \mathcal{S} de sommets de $V \setminus \mathcal{N}$.
 Tant que nous n'avons pas effectué T tirages, on pioche $x \in \mathcal{S}$, et on réutilise la stratégie avec τ connu en l'approximant par $\hat{\tau}$, le tout sur les communautés \hat{C}_0 et \hat{C}_1 :
 - On tire un sommet dans chaque \hat{C} .
 - Tant que ni \hat{C}_0 ni \hat{C}_1 n'a été intégralement visitée :
 - Si les $\hat{\mu}$ sont tous deux inférieurs à $\hat{\tau}$, on tire un sommet supplémentaire dans chaque \hat{C} .
 - Sinon, on tire un sommet de la communauté associée au $\hat{\mu}$ le plus élevé.
- **On ajoute x à la communauté estimée dans laquelle on a effectué le dernier tirage.**

2.3.1 • CONJECTURES SUR CET ALGORITHME

On pense que cet algorithme, très similaire au premier à ceci près que le noyau croît après chaque nouvelle classification, donne de meilleurs résultats asymptotiques que l'algorithme précédent.

Plus précisément, on conjecture que cette méthode donne un regret en $O\left(\frac{p}{\Delta}\sqrt{T}\right)$, qui est la borne optimale pour ce modèle, *cf.* [2].

En effet, nos simulations suggèrent clairement que la deuxième stratégie proposée est la meilleure :



Donnons maintenant quelques arguments en faveur de ce résultat.

2.3.2 • RETOUR SUR LE CHOIX DE n

On a vu dans la partie précédente que le regret respectait une inégalité de la forme :

$$R_T \leq \Delta \left(\binom{n}{2} + \alpha \frac{p}{\Delta^2} |\mathcal{S}| \right) + \varepsilon(T)$$

avec $\varepsilon(T) \xrightarrow{T \rightarrow +\infty} 0$ et α une constante.

Le fait d'ajouter les sommets à nos communautés pourrait permettre de réduire drastiquement les valeurs de n et $|\mathcal{S}|$ optimales.

Dans l'algorithme précédent, la meilleure borne pour $|\mathcal{S}|$ était de la forme :

$$|\mathcal{S}| \leq \alpha \frac{T}{n}$$

avec α une constante, ce qui aboutissait à une valeur de n de l'ordre de $T^{1/3}$.

Nous pensons qu'ajouter des sommets au fur et à mesure permet d'arriver à une borne de la forme :

$$|\mathcal{S}| \leq \alpha \sqrt{T}$$

Pour peu que $n \geq T^\varepsilon$ pour un certain $\varepsilon > 0$, on aurait alors :

$$R_T \leq \Delta \left(\binom{n}{2} + \alpha \frac{p}{\Delta^2} \sqrt{T} \right) + \varepsilon(T)$$

En prenant $n = T^{1/4}$, on trouverait bien :

$$R_T = O \left(\frac{p}{\Delta} \sqrt{T} \right)$$

2.3.3 • MAJORATION DE $|\mathcal{S}|$

On s'intéresse ici à \mathcal{S}^* , qu'on définit comme égal à \mathcal{S} privé d'un éventuel dernier sommet problématique : rien n'indique que notre algorithme épuisera une classification au T^e tour.

En supposant que notre deuxième algorithme classe **parfaitement** tous les sommets de \mathcal{S}^* , on se convainc rapidement que chaque sommet de $\mathcal{S}^* \cap C_i$ finit par voir toutes ses arêtes avec les sommets de \widehat{C}_i et $\mathcal{S}^* \cap C_i$ tirées.

On aurait donc au moins

$$\binom{|\mathcal{S}^* \cap C_i|}{2} + |\widehat{C}_i| \cdot (|\mathcal{S}^* \cap C_i|)$$

tirages pour chaque communauté.

Cela donne au moins $\frac{|\mathcal{S}^*|^2}{4}$ tirages, toutes communautés confondues.

$|\mathcal{S}| = \lceil 2\sqrt{T} \rceil + 1$ sommets suffisent donc pour atteindre les T tirages.

2.3.4 • RÉSULTATS À PROUVER

Il reste à prouver que la catégorisation faite par ce deuxième algorithme est suffisamment bonne pour garder une proportion suffisamment faible de sommets mal catégorisés dans les communautés estimées au fil des ajouts de sommets. C'est à dire qu'à tout instant et pour $i \in \{0, 1\}$ on ait l'existence d' $\varepsilon > 0$ tel que :

$$\frac{|\hat{C}_i \cap C_{1-i}|}{|\hat{C}_i|} \leq \varepsilon < \frac{1}{2} - \varepsilon_\tau$$

Cette inégalité serait suffisante à ce que la quantité

$$1 + \frac{6p}{\Delta^2 \left(\frac{1}{2} - \varepsilon - \varepsilon_\tau \right)^2}$$

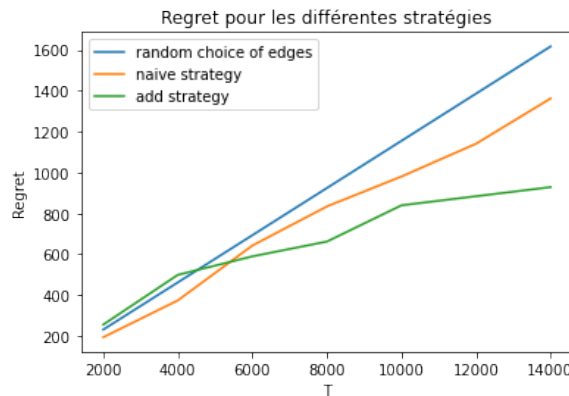
majoré l'espérance du nombre de mauvais tirages pour un sommet donné, ce qui nous permettrait de conclure.

À ce stade, nous n'avons trouvé aucun moyen piste pour montrer que la qualité des communautés ne se détériore pas trop pour des valeurs de T plus grande, ce que nos simulations semblent pourtant suggérer.

2.3.5 • REMARQUES D'ORDRE PRATIQUE

Les détails techniques de notre implémentation sont consultables dans notre notebook Jupyter.

Voici les courbes obtenues pour nos données issues de Reddit :



Les données corroborent nos conclusions et aboutissent sensiblement aux mêmes résultats que nos simulations.

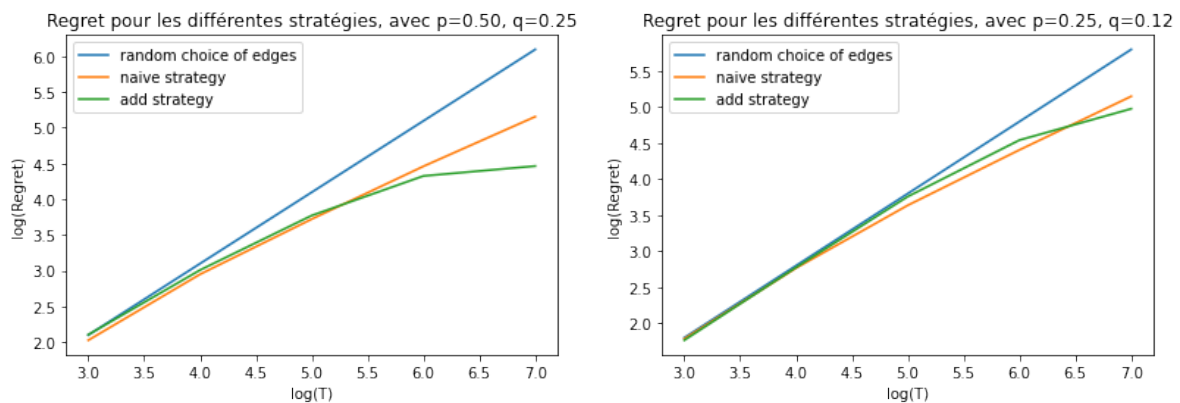
Cependant plusieurs remarques peuvent être faites, en sachant que nos valeurs de p et q sont à peu près similaires à celles choisies pour nos simulations : $p \simeq 0,32$ et $q \simeq 0,07$.

La première est que la deuxième stratégie n'a que très peu d'intérêt, voire est contre-productive pour des valeurs de T trop faibles.

La deuxième est que nos données n'ont *a priori* aucune raison de coller à notre modèle. En effet certains utilisateurs plus actifs sont à l'origine d'un plus grand nombre d'arêtes. Les valeurs estimées de p et q ne sont donc que des moyennes, et la variabilité dans le nombre de connexions aux autres sommets de la communauté est plus grande que dans nos simulations.

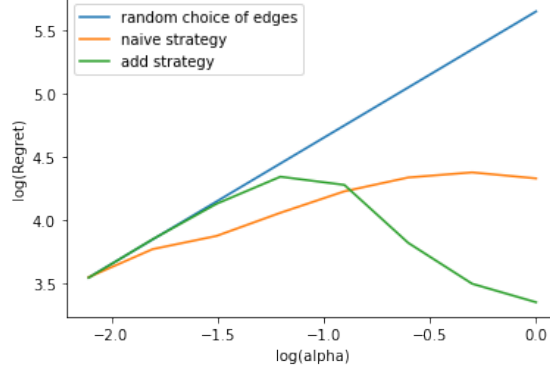
La troisième est que malgré tout ce que l'on a dit, on parvient à diminuer significativement notre regret par rapport à une sélection aléatoire d'arêtes.

On peut ajouter les résultats suivants :



Ainsi, le fait de réduire nos paramètres p et q d'un même facteur pousse vers le premier régime, celui où nos stratégies ne font guère mieux que le hasard.

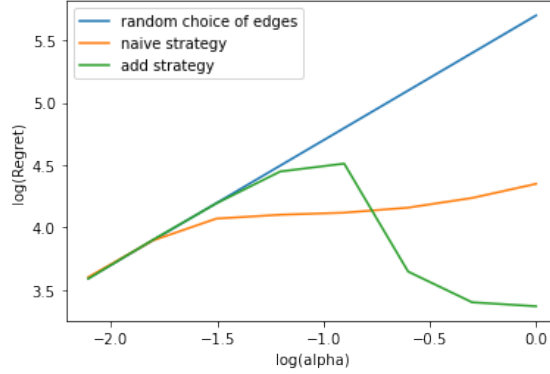
On peut encore mieux observer ce phénomène dans le graphique suivant :

Regret pour les différentes stratégies, avec $T=1e6$, $p=0.95*\alpha$, $q=0.05*\alpha$ 

On choisit ici d'étudier, pour $T = 10^6$, ce qu'il se passe pour des paramètres $p = \alpha p_0$ et $q = \alpha q_0$, avec $p_0 = 0,95$, $q_0 = 0,05$ et $\alpha \in [0, 1]$.

Il nous permet notamment de voir que pour un T de l'ordre du million, si nos valeurs de p et q sont trop faibles, il vaut parfois mieux utiliser notre premier algorithme, et ceux même si l'écart relatif entre p et q est important (il est ici constant et égal à 20).

Similairement, on peut fixer $p = p_0$ et poser $q = (1 - \alpha)p_0$, ce qui revient à faire varier Δ . On trouve alors, pour $p_0 = 0,95$:

Regret pour les différentes stratégies, avec $T=1e6$, $p=0.95$, $q=0.95*(1-\alpha)$ 

On retrouve une dynamique très similaire, avec là encore une performance meilleure pour le premier algorithme lorsque p et q sont trop proches compte tenu de la valeur de T .

Que cela soit à p fixé ou non, cette meilleure performance du premier algorithme était attendue : pour $T = 10^6$, le noyau sur lequel se base le deuxième algorithme comporte $\sqrt[4]{10^6} < 32$ éléments. Si p et q sont trop proches, notre marge d'erreur sur la classification de ce noyau sera trop grande, et les erreurs commises se répercuteront sur toutes les classifications opérées durant la deuxième phase.

RÉFÉRENCES

- [1] R. J. Serfling. Probability Inequalities for the Sum in Sampling without Replacement. *The Annals of Statistics*, 2(1) : 39 – 48, 1974.
- [2] Christophe Giraud, Yann Issartel, Luc Lehericy, and Matthieu Lerasle. Pair-matching : Links prediction with adaptive queries, 2020.