



Journées mathématiques X-UPS

Année 2016

Arbres et marches aléatoires

Grégory MIERMONT

Probabilités sur le graphe complet : l'exemple des arbres couvrants uniforme et minimal

Journées mathématiques X-UPS (2016), p. 59-102.

<https://doi.org/10.5802/xups.2016-02>

© Les auteurs, 2016.



Cet article est mis à disposition selon les termes de la licence

LICENCE INTERNATIONALE D'ATTRIBUTION CREATIVE COMMONS BY 4.0.

<https://creativecommons.org/licenses/by/4.0/>

Les Éditions de l'École polytechnique
Route de Saclay
F-91128 PALAISEAU CEDEX
<https://www.editions.polytechnique.fr>

Centre de mathématiques Laurent Schwartz
CMLS, École polytechnique, CNRS,
Institut polytechnique de Paris
F-91128 PALAISEAU CEDEX
<https://portail.polytechnique.edu/cmls/>



Publication membre du

Centre Mersenne pour l'édition scientifique ouverte

www.centre-mersenne.org

PROBABILITÉS SUR LE GRAPHE COMPLET : L'EXEMPLE DES ARBRES COUVRANTS UNIFORME ET MINIMAL

par

Grégory Miermont

Résumé. Ce texte propose quelques exemples d'analyse de grandes structures combinatoires aléatoires, que l'on peut définir naturellement en termes de modèles simples d'arbres couvrants sur le graphe complet.

Table des matières

1. Introduction.....	60
1.1. Contexte.....	60
1.2. Terminologie et notations.....	61
2. L'arbre couvrant uniforme.....	63
2.1. Dénombrement.....	64
2.2. Simulation.....	64
2.3. Une autre interprétation.....	72
2.4. Géométrie asymptotique : l'arbre continu aléatoire	73
2.5. Perspectives sur le CRT.....	79
<i>Intermezzo</i> : le processus des forêts coalescentes de Pitman	80
3. L'arbre couvrant minimal.....	84
3.1. L'algorithme de Kruskal, et applications.....	85
3.2. Le « théorème $\zeta(3)$ » de Frieze.....	88
3.3. La convergence locale.....	90
3.4. Le graphe complet pondéré vu comme un arbre...	92
3.5. L'arbre minimal vu comme une forêt couvrante du PWIT.....	95
3.6. Géométrie asymptotique.....	99
Références.....	100

1. Introduction

1.1. Contexte

Le thème développé dans ce texte se situe à l'interface entre différents aspects des probabilités : physique statistique, combinatoire, optimisation. Le thème de la physique statistique est de décrire l'état d'un « grand système » à l'aide d'une description microscopique, au niveau particulaire : souvent, pour simplifier ou pour des raisons physiques, on considère que ces espaces de configurations sont définis sur les sites d'un graphe fini, par exemple un sous-graphe de \mathbb{Z}^d . Quelques exemples de modèles de référence sont

- La percolation, où chaque arête du graphe est conservée avec une probabilité donnée, indépendamment entre les arêtes
- Le modèle de Lenz-Ising du ferromagnétisme, où chaque sommet est muni d'un signe $+$ ou $-$, et où chaque configuration apparaît avec une probabilité dépendant d'une énergie définie en termes du nombre de sommets adjacents de signes opposés
- L'arbre couvrant uniforme, qui est le modèle que nous allons étudier plus en détail
- La marche aléatoire auto-évitante uniforme, qui est un modèle naturel décrivant le repliement d'un polymère, traité dans les exposés de Vincent Beffara.

Pour une introduction accessible à certains de ces sujets, on recommande notamment la lecture des excellents articles de Raphaël Cerf, Hugo Duminil-Copin et Marie Thérêt sur le site web *Images des mathématiques*. Par ailleurs, en optimisation combinatoire, on cherche à minimiser une fonction aléatoire de l'espace des configurations (souvent interprétée comme une énergie), avec l'idée de décrire un système physique dans un état d'équilibre. On s'intéressera ici encore au cas où l'espace des configurations est l'ensemble des arbres couvrants d'un graphe donné.

Nous allons considérer de tels systèmes définis sur le graphe complet, où chaque sommet est connecté à tous les autres. Cette « géométrie » est beaucoup plus simple à décrire que celle du réseau \mathbb{Z}^d , par exemple par un comptage entièrement explicite des configurations

possibles, d'où le lien avec la combinatoire. Mais il est instructif d'étudier les modèles ci-dessus dans ce cas simple : à la fois parce qu'ils peuvent donner une intuition précieuse sur les modèles sur réseau, mais aussi parce que leur simplicité fait surgir des objets fondamentaux que l'on retrouve en de nombreuses autres occasions en probabilités. Les deux objets-clés qui nous intéresseront dans ces notes ont été introduits par David Aldous [Ald91, Ald92]. Ce sont deux modèles d'arbres aléatoires « continus » de natures très différentes. Le premier (l'arbre continu aléatoire brownien) est également un acteur-clé du texte d'Igor Kortchemski. Le second est une intrigante structure arborescente où tout sommet est à la fois de degré infini, et n'a pourtant qu'un nombre fini de voisins à distance bornée. Nous suivrons d'assez près l'article [Ald91], ainsi que l'article de revue d'Aldous et Steele [AS04].

1.2. Terminologie et notations

On emploiera les notations $\mathbb{N} = \{1, 2, \dots\}$ et $\mathbb{Z}_+ = \mathbb{N} \cup \{0\}$. Pour $n \in \mathbb{N}$, on notera $[n] = \{1, 2, \dots, n\}$.

Graphes. Rappelons qu'un *graphe* (simple) est un couple $G = (V, E)$, où $V = V(G)$ est un ensemble de *sommets*, et $E = E(G)$ est un ensemble de paires $\{x, y\}$ avec $x, y \in V$ et $x \neq y$. On notera souvent $|G| = \text{Card}(G)$ le cardinal de l'ensemble des *arêtes* de G .

On dit que les sommets x, y sont *adjacents* si $e = \{x, y\} \in E$, ce que l'on note $x \sim y$ pour simplifier, et l'on dit que x, y sont les extrémités de l'arête e . Le graphe est *orienté* si l'on distingue un ordre $\vec{e} = (x, y)$ pour chaque arête $e = \{x, y\} \in E$. On dit alors que x est l'*origine*, et y la *cible* de \vec{e} .

Un chemin dans un graphe $G = (V, E)$ est une suite finie (x_1, x_2, \dots, x_k) de sommets avec

$$x_i \sim x_{i+1} \quad \text{pour tout } i \in \{1, 2, \dots, k-1\},$$

et on dit que $k-1$ est la *longueur* de ce chemin. On dit que c'est une *chaîne* si de plus tous les x_i sont deux à deux distincts, et un *cycle* si $x_1 = x_n$ et si x_1, \dots, x_{n-1} sont deux à deux distincts. Le graphe G

est *connexe* si pour tout $x, y \in V$, il existe une chaîne commençant en x et finissant un y .

Un *arbre* est un graphe $G = (V, E)$ connexe, sans cycle. Si V est fini, cela équivaut au fait que G est connexe et a $\text{Card}(V) - 1$ arêtes, ou encore, que G est sans cycle et a $\text{Card}(V) - 1$ arêtes. Ainsi, $n - 1$ est le nombre minimal (resp. maximal) d'arêtes permettant de connecter n sommets (resp. permettant de ne pas créer de cycle sur n sommets).

Un *graphe enraciné* est un couple (G, v) où $G = (V, E)$ est un graphe et $v \in V$. Un *graphe marqué* est un couple (G, f) où $f : E \rightarrow \Xi$ est une fonction positive définie sur les arêtes, à valeurs dans un ensemble de marques Ξ . Si $\Xi = (0, \infty)$, on parlera de *graphe pondéré*. Un sous-graphe *couvrant* de G est un graphe $G' = (V, E')$ avec $E' \subset E$: on peut également le voir comme le graphe marqué (V, E', f) avec $f(e) = \mathbf{1}_{\{e \in E'\}}$. Si G' est un sous-graphe couvrant qui est un arbre, on dit que G' est un arbre couvrant.

On notera $K_n = ([n], E(K_n))$ le graphe complet à n sommets étiquetés $[n] = \{1, 2, \dots, n\}$, et avec $E(K_n) = \{\{i, j\} : i, j \in [n], i \neq j\}$.

Notions probabilistes. Les variables aléatoires que nous considérerons seront toujours supposées définies sur un espace de probabilités $(\Omega, \mathcal{F}, \mathbb{P})$.

Soit $(X_n, n \geq 1)$, X des variables aléatoires à valeurs dans un espace métrique (M, d) . On dit que X_n converge *en loi* vers X si pour toute fonction $f : M \rightarrow \mathbb{R}$ continue bornée, on a

$$\mathbb{E}[f(X_n)] \xrightarrow{n \rightarrow \infty} \mathbb{E}[f(X)].$$

Dans le cas où $M = \mathbb{R}^d$ normé, il suffit en fait de le vérifier pour les fonctions f continues à support compact. Voici quelques faits sur la convergence en loi qui nous serviront au cours du texte.

(1) Par un théorème de Skorokhod, la convergence en loi de X_n vers X signifie que l'on peut trouver, sur un espace de probabilités éventuellement différent, des variables aléatoires $(X'_n, n \geq 1)$, X' telles que

- pour tout $n \geq 1$, X'_n a même loi que X_n
- X' a même loi que X

• l'événement $\lim X'_n = X'$ lorsque $n \rightarrow \infty$ a probabilité 1 (il est presque sûr).

(2) Si X_n est à valeurs dans \mathbb{R} , il y a équivalence entre la convergence en loi de X_n vers X et le fait que $\mathbb{P}(X_n > r)$ converge vers $\mathbb{P}(X > r)$, pour tout r tel que $\mathbb{P}(X = r) = 0$.

(3) Si X_n est à valeurs dans \mathbb{Z}^d pour un entier $d \geq 1$, et si l'on a que pour une suite $a_n \rightarrow \infty$

$$\mathbb{P}(X_n = x) = \frac{1}{a_n^d} g(x/a_n) (1 + o(1)),$$

ce uniformément pour tout $x = O(a_n)$, où $g : \mathbb{R}^d \rightarrow \mathbb{R}_+$ est une densité de probabilités (c'est-à-dire que $\int g = 1$), alors on a que X_n/a_n converge en loi vers X , variable aléatoire dont la loi est donnée par $g(x)dx$ sur \mathbb{R}^d .

2. L'arbre couvrant uniforme

Soit $G = (V, E)$ un graphe fini (simple, non orienté) et connexe. Un arbre couvrant est un sous-graphe $T = (V, E')$, $E' \subset E$, qui est un arbre. On note \mathbf{A}_G l'ensemble des arbres couvrants de G , qui est évidemment fini. On dit qu'une variable aléatoire A sur un espace de probabilités $(\omega, \mathcal{F}, \mathbb{P})$, et à valeurs dans \mathbf{A}_G , est un arbre couvrant uniforme de G si

$$\mathbb{P}(A = \mathbf{a}) = \frac{1}{\text{Card}(\mathbf{A}_G)}, \quad \mathbf{a} \in \mathbf{A}_G.$$

Dans le cas où $G = K_n$, l'ensemble $\mathbf{A}_n = \mathbf{A}_{K_n}$ est tout simplement l'ensemble des arbres sur $[n]$. On peut se demander pourquoi l'on insiste sur le fait que ces arbres « couvrent » K_n , la raison étant que l'ensemble des arbres couvrants d'un graphe est extrêmement étudié. Pour commencer, on peut simplement se demander combien il y en a.

Théorème 2.1 (Cayley). *Le cardinal de \mathbf{A}_n est n^{n-2} .*

Il existe de très nombreuses preuves de ce résultat, et nous en rencontrerons trois chemin faisant, une partielle et deux complètes.

Pour $v_0 \in V$ on notera également $\mathbf{A}_G(v_0)$ l'ensemble des arbres couvrants (\mathbf{a}, v_0) de G enracinés en v_0 , et dans ce cas, on conviendra

que les arêtes $\{u, v\}$ de \mathbf{a} sont systématiquement munies de l'orientation (u, v) telle que $d_{\mathbf{a}}(u, v_0) = d_{\mathbf{a}}(v, v_0) + 1$, où $d_{\mathbf{a}}(v, w)$ est la distance de graphe sur \mathbf{a} entre v et w . Une feuille de \mathbf{a} est par définition un sommet vers lequel aucune arête ne pointe, dans cette orientation canonique.

2.1. Dénombrement

On peut déterminer le cardinal de \mathbf{A}_G en termes du déterminant du laplacien du graphe : soit M la matrice d'adjacence de $G = (V, E)$, qui est la matrice définie par

$$M(v, v') = \mathbf{1}_{\{v \sim v'\}}, \quad v, v' \in V,$$

où l'on note $v \sim v'$ si $\{v, v'\} \in E$, et notons $D(v, v') = \deg(v)\mathbf{1}_{\{v=v'\}}$ où $\deg(v) = \text{Card}\{v' \in V : v' \sim v\}$ est le degré de v . Enfin, on note $\Delta^G = D - M$.

Clairement, Δ^G n'est pas une matrice inversible puisque $\Delta^G \mathbf{1} = 0$, où $\mathbf{1}$ est le vecteur de \mathbb{R}^V dont toutes les composantes sont égales à 1. En revanche, on peut montrer que si G est connexe, alors la matrice $\Delta_{v_0}^G$ obtenue en retirant la ligne et la colonne correspondant à $v_0 \in V$ est inversible pour tout v_0 , et de plus

Théorème 2.2 (Kirchhoff). *On a que $|\det(\Delta_{v_0}^G)| = \text{Card}(\mathbf{A}_G)$ pour tout $v_0 \in V$.*

Nous ne donnons pas la preuve de ce théorème ici. Une application directe au cas de $G = K_n$ et $v_0 = 1$ donne une première preuve du théorème de Cayley ci-dessus, par le calcul du déterminant de la matrice circulante $(n-1) \times (n-1)$:

$$\Delta_1^{K_n} = \begin{pmatrix} n-1 & -1 & \dots & -1 \\ -1 & n-1 & \dots & -1 \\ \vdots & \vdots & \ddots & \vdots \\ -1 & -1 & \dots & n-1 \end{pmatrix}.$$

2.2. Simulation

Nous cherchons à présent à décrire la géométrie d'un arbre couvrant aléatoire de K_n . Par exemple, quel est le diamètre typique

d'un tel arbre, en fonction de n , et comment ce diamètre évolue-t-il lorsque $n \rightarrow \infty$? Un problème naturel similaire consiste à étudier le comportement asymptotique d'une marche aléatoire

$$S_n = X_1 + X_2 + \cdots + X_n$$

où les variables aléatoires X_1, X_2, \dots sont indépendantes et de même loi, par exemple $\mathbb{P}(X_1 = 1) = 1 - P(X_1 = -1) = 1/2$. En effet, la loi forte des grands nombres stipule que $S_n/n \rightarrow 0$ presque sûrement lorsque $n \rightarrow \infty$, tandis que le théorème central limite donne

$$\mathbb{P}\left(a \leq \frac{S_n}{\sqrt{n}} \leq b\right) \xrightarrow{n \rightarrow \infty} \int_a^b \sqrt{\frac{2}{\pi}} \exp(-2x^2) dx$$

pour tout $a < b$. Ainsi, la marche aléatoire S_n est « typiquement » de l'ordre de \sqrt{n} au sens où la probabilité que S_n n'est pas dans $[-K\sqrt{n}, K\sqrt{n}]$ tend vers 0 lorsque $K \rightarrow \infty$ uniformément en n .

2.2.1. L'algorithme d'Aldous-Broder. Pour pouvoir répondre à ce genre de question pour un arbre uniforme, il ne suffit évidemment pas de connaître le nombre de possibilités, comme il serait insuffisant dans le cas de la marche aléatoire de savoir que S_n est déterminée par les 2^n valeurs que peut prendre (X_1, X_2, \dots, X_n) . On est donc amené à essayer de décrire la structure d'un arbre couvrant uniforme, et pour cela on va avoir recours à un algorithme permettant d'en simuler un. Ce dernier a été proposé par Aldous [Ald90] et Broder au début des années 1990.

Théorème 2.3. *Soit $v_0 \in V$ fixé. On se donne une marche aléatoire $X_0 = v_0, X_1, X_2, \dots$ au plus proche voisin dans G , issue de v_0 , c'est-à-dire que pour tout $x_0 = v_0, x_1, \dots, x_i$ dans V ,*

$$\mathbb{P}(X_{i+1} = y | X_0 = x_0, X_1 = x_1, \dots, X_i = x_i) = \frac{1}{\deg(x_i)} \mathbf{1}_{\{y \sim x_i\}},$$

dès lors que l'événement par lequel on conditionne est de probabilité non nulle. Soit $T_v = \inf\{i \geq 0 : X_i = v\}$ le premier temps d'atteinte de $v \in V$ par cette marche. Alors le graphe aléatoire $A = (V, \{\{X_{T_v-1}, X_{T_v}\} : v \in V \setminus \{v_0\}\})$, vu comme enraciné en v_0 , est un arbre couvrant uniforme parmi l'ensemble $\mathbf{A}_G(v_0)$ des arbres couvrants de G enracinés en v_0 .

On notera que les temps T_v sont tous finis presque sûrement, par un résultat classique sur les chaînes de Markov, dont nous n'aurons pas besoin explicitement dans ces notes et que nous ne discuterons pas davantage. Il faut prendre garde dans l'énoncé précédent qu'il n'est pas vrai en général que le graphe A « dé-raciné » est uniforme dans \mathbf{A}_G . Cela devient néanmoins vrai si l'on a aussi choisi le sommet v_0 aléatoirement, selon la loi $\pi(v) = \deg(v) / \sum_{w \in V} \deg(w)$ sur V .

Donnons un point de vue un peu différent sur ce résultat, qui met en évidence son caractère algorithmique.

Algorithme d'Aldous-Broder.

- (1) Initialiser $S = \{v_0\}$, $Z = \emptyset$, $I = 0$, $y = v_0$.
- (2) Tant que $S \neq V$, faire
 - $I \leftarrow I + 1$,
 - si $X_I \notin S$ alors $Z \leftarrow Z \cup \{\{y, X_I\}\}$,
 - $S \leftarrow S \cup \{X_I\}$
 - $y \leftarrow X_I$
- (3) Retourner Z .

Autrement dit, on construit un sous-graphe de G en traçant les arêtes visitées par la marche aléatoire les unes après les autres, sauf aux moments où l'on retourne sur un sommet déjà visité auparavant : en de tels instants, on se place en ce sommet, mais sans ajouter l'arête qui y mène. L'algorithme termine pour la même raison que précédemment, au premier instant C où $\{v_0, X_1, X_2, \dots, X_C\} = V$ (ce temps C est appelé *temps de couverture* issu de v_0).

Notons $R_0 = 0$, et par récurrence, pour tout $i \geq 0$, soit

$$R_{i+1} = \inf\{j > i : X_j \in \{X_0, X_1, \dots, X_{j-1}\}\},$$

c'est-à-dire que R_1, R_2, \dots sont les instants successifs où la suite X_0, X_1, \dots répète une valeur déjà prise par le passé. L'arbre A de l'algorithme peut alors s'interpréter comme le graphe dont les arêtes sont

$$\{\{X_{j-1}, X_j\} : j \geq 1, j \notin \{R_1, R_2, \dots\}\}.$$

2.2.2. Le cas du graphe complet. Nous allons donner une preuve du théorème 2.3 *uniquement dans le cas où $G = K_n$ est le graphe complet*. Dans ce cas précis, la marche aléatoire (X_1, X_2, \dots) est particulièrement simple : à l'étape i , on choisit un élément de $[n]$ distinct de X_i uniformément au hasard. En fait, pour simplifier les choses encore un peu, nous allons supposer que X_1, X_2, \dots est une suite de variables aléatoires i.i.d. uniformes dans $[n]$: cela revient à dire que la marche aléatoire reste sur place avec probabilité $1/n$ à chaque étape. Clairement, cela ne modifie pas la nature de l'algorithme d'Aldous-Broder.

Par ailleurs, pour fixer les idées, on prendra $X_0 = 1$ dans l'énoncé qui suit. Néanmoins, clairement, on peut oublier l'enracinement par symétrie. Suivant Camarri-Pitman [CP00], on peut donner une description exacte du déroulé de l'algorithme, au sens suivant.

Lemme 2.4. *Soit $m \geq 1$ fixé, et $y_1, y_2, \dots, y_m, x \in [n]$. On se donne un arbre $\mathbf{a} = (V(\mathbf{a}), E(\mathbf{a}))$ enraciné en 1 sur un sous-ensemble $V(\mathbf{a})$ de $[n]$, dont les feuilles sont incluses dans $\{y_1, \dots, y_m\}$. On note enfin $A[m]$ le sous-arbre de A obtenu en ne gardant que les arêtes $\{X_{j-1}, X_j\}$ avec $j \in \{1, 2, \dots, R_m - 1\} \setminus \{R_1, R_2, \dots, R_{m-1}\}$. Alors si $|\mathbf{a}|$ est le nombre d'arêtes de \mathbf{a} , et si x est un sommet de \mathbf{a} ,*

$$\mathbb{P}(X_{R_i-1} = y_i, 1 \leq i \leq m; X_{R_m} = x; A[m] = \mathbf{a}) = \frac{1}{n^{m+|\mathbf{a}|}}.$$

Démonstration. Notons que la contrainte que toute feuille de \mathbf{a} soit de la forme y_i provient de la construction de l'algorithme : en effet, si i n'est pas un instant de la forme $R_j - 1$ pour un $j \geq 1$, cela signifie que $i + 1$ est le premier instant où le sommet X_{i+1} est visité, et l'arête $\{X_i, X_{i+1}\}$ est dans l'arbre construit par l'algorithme, avec l'orientation (X_{i+1}, X_i) . Donc X_i n'est pas une feuille de A .

Notons aussi que l'on a $R_m = m + |A[m]|$ puisque $A[m]$ est constitué de $R_m - m$ arêtes (une arête pour chaque itération de l'algorithme, sauf lorsque cela implique une répétition).

On montre alors le résultat par récurrence sur m . Pour $m = 1$, l'arbre $A[1]$ est la chaîne $1 = X_0, X_1, \dots, X_{R_1-1}$. De ce fait, si \mathbf{a} est une chaîne $1 = x_0, x_1, \dots, x_k$, si $y_1 = x_k$ est l'unique feuille de \mathbf{a} (enraciné en 1), et si $x \in \{x_0, x_1, \dots, x_k\}$, l'événement $\{X_{R_1-1} = y_1, X_{R_1} =$

$x, A[1] = \mathbf{a}$ est l'événement $\{X_1 = x_1, \dots, X_k = x_k, X_{k+1} = x\}$, et cela a probabilité $1/n^{k+1} = 1/n^{1+|\mathbf{a}|}$.

Supposons le résultat connu au rang m . Notons alors que $A[m+1]$ s'obtient en ajoutant la chaîne $\{X_{R_m}, X_{R_m+1}, \dots, X_{R_{m+1}-1}\}$ au graphe $A[m]$. Soit alors \mathbf{a} un arbre enraciné en 1 dont les feuilles sont incluses dans un ensemble $\{y_1, \dots, y_{m+1}\}$, et soit \mathbf{a}' le sous-arbre engendré par la racine et les sommets y_1, \dots, y_m : il s'agit de la réunion des chaînes orientées de y_1, \dots, y_m vers la racine 1 dans \mathbf{a} . La chaîne orientée de y_{m+1} vers 1 rencontre \mathbf{a}' pour la première fois en un sommet x' , de sorte que \mathbf{a} est la réunion de \mathbf{a}' avec une chaîne, disons $\{x' = x_0, x_1, \dots, x_k\}$, avec $k = |\mathbf{a}| - |\mathbf{a}'|$. L'événement $\{X_{R_i-1} = y_i, 1 \leq i \leq m+1; A[m+1] = \mathbf{a}, X_{R_{m+1}} = x\}$ peut alors se réécrire

$$\{X_{R_i-1} = y_i, 1 \leq i \leq m; A_m = \mathbf{a}', X_{R_m} = x'\} \\ \cap \{X_{R_m+1} = x_1, \dots, X_{R_m+k} = x_k, X_{R_m+k+1} = x\},$$

mais sur le premier événement, on a que $R_m = m + |\mathbf{a}|$, et on voit que le premier événement ne dépend que des variables aléatoires

$$X_1, X_2, \dots, X_{m+|\mathbf{a}|},$$

tandis que le second dépend de

$$X_{m+1+|\mathbf{a}|}, X_{m+2+|\mathbf{a}|}, \dots$$

Par indépendance, on conclut que

$$\begin{aligned} \mathbb{P}(X_{R_i-1} = y_i, 1 \leq i \leq m+1; A[m+1] = \mathbf{a}, X_{R_{m+1}} = x) \\ &= \frac{1}{n^{|\mathbf{a}|-|\mathbf{a}'|+1}} \mathbb{P}(X_{R_i-1} = y_i, 1 \leq i \leq m; A_m = \mathbf{a}', X_{R_m} = x') \\ &= \frac{1}{n^{|\mathbf{a}|-|\mathbf{a}'|+1}} \cdot \frac{1}{n^{m+|\mathbf{a}'|}} = \frac{1}{n^{m+1+|\mathbf{a}|}} \end{aligned}$$

où l'on a utilisé l'hypothèse de récurrence à l'avant-dernière étape. D'où le résultat. \square

De ce résultat, on déduit le théorème 2.3 dans le cas particulier $G = K_n$.

Corollaire 2.5. *La probabilité que l'algorithme d'Aldous-Broder sur le graphe complet donne un arbre couvrant \mathbf{a} donné est*

$$\mathbb{P}(A = \mathbf{a}) = \frac{1}{n^{n-2}}.$$

En particulier, A est uniforme dans \mathbf{A}_n et $\text{Card}(\mathbf{A}_n) = n^{n-2}$. De plus, la suite $X_{R_1-1}, X_{R_2-1}, \dots$ est i.i.d. uniforme sur $\{1, 2, \dots, n\}$, et indépendante de A

Démonstration. Soit \mathbf{a} un arbre couvrant de K_n , et $m \geq n$. Si l'on somme la formule du lemme 2.4 sur toutes les familles y_1, \dots, y_m telles que $\{y_1, \dots, y_m\} = [n]$ (on notera $\mathbf{y} = (y_1, \dots, y_m) \in C_{m,n}$), alors en notant que $A[m] = A$ sur l'événement où $(X_{R_i-1}, 1 \leq i \leq m) \in C_{m,n}$,

$$\begin{aligned} & \mathbb{P}((X_{R_i-1}, 1 \leq i \leq m) \in C_{m,n}; A = \mathbf{a}) \\ &= \sum_{\substack{\mathbf{y} \in C_{m,n} \\ x \in [n]}} \mathbb{P}(X_{R_i-1} = y_i, 1 \leq i \leq m; X_{R_m} = x; A[m] = \mathbf{a}) \\ &= \frac{n \text{Card } C_{m,n}}{n^{m+1}} \cdot \frac{1}{n^{|\mathbf{a}|-1}} = \frac{\mathbb{P}((X_1, \dots, X_m) \in C_{m,n})}{n^{n-2}}. \end{aligned}$$

Ici, on a utilisé le fait que $|\mathbf{a}| = n - 1$, qui vient du fait général qu'un arbre a toujours un sommet de plus que d'arêtes. Si l'on fait tendre m vers l'infini, chacun des événements $\{(X_{R_i-1}, 1 \leq i \leq m) \in C_{m,n}\}$ et $(X_i, 1 \leq i \leq m) \in C_{m,n}$ a une probabilité qui tend vers 1, n étant fixé. On conclut donc par un passage à la limite que l'on a la formule attendue.

La fin de l'énoncé est obtenu de façon similaire : cette fois on choisit $k \geq 1$ et des entiers $y_1, \dots, y_k \in [n]$ et on écrit, pour tout $m \geq n$,

$$\begin{aligned} & \mathbb{P}((X_{R_i-1}, k < i \leq m+k) \in C_{m,n}; X_{R_i-1} = y_i, 1 \leq i \leq k; A = \mathbf{a}) \\ &= \sum_{\mathbf{y}' \in C_{m,n}} \mathbb{P}(X_{R_i-1} = y_i, 1 \leq i \leq k; X_{R_{k+i}-1} = y'_i, 1 \leq i \leq m; A = \mathbf{a}) \\ &= \frac{\mathbb{P}((X_1, \dots, X_m) \in C_{m,n})}{n^k \cdot n^{n-2}}, \end{aligned}$$

enfin, on prend la limite $m \rightarrow \infty$ et on conclut. \square

2.2.3. Efficacité. Combien de temps faut-il à l'algorithme pour terminer ? Dans le cas du graphe complet, le temps de couverture est (à une soustraction de 1 près, puisque l'on a $X_0 = 0$)

$$C_n = \inf\{k \geq 1 : \{X_1, \dots, X_k\} = [n]\}.$$

On montre que ce temps C_n a même loi qu'une somme de variables aléatoires indépendantes de lois géométriques, respectivement de paramètres $1, (n-1)/n, (n-2)/n, \dots, 1/n$:

$$C_n = \sum_{i=1}^n \tau_i$$

où

$$\mathbb{P}(\tau_i > k) = \left(\frac{i-1}{n}\right)^k, \quad i \in [n], k \geq 0.$$

Ici, $\tau_i - 1$ représente le nombre de fois où la suite X_1, X_2, \dots a répété une des $i-1$ premières valeurs déjà prises avant de découvrir pour la première fois une i -ième valeur. Une application d'une inégalité de Bienaymé-Chebychev montre que l'on a que $C_n/n \ln n$ converge vers 1 en probabilité, c'est-à-dire que pour tout $\varepsilon > 0$,

$$\mathbb{P}\left(\left|\frac{C_n}{n \ln(n)} - 1\right| > \varepsilon\right) \xrightarrow{n \rightarrow \infty} 0.$$

On parle du problème du *collectionneur de coupons* : C_n est le temps que doit mettre un collectionneur pour obtenir une collection complète de n images, s'il se procure à chaque instant une image choisie uniformément au hasard dans la collection.

On peut également se demander à quelle vitesse la probabilité ci-dessus tend vers 0. Pour ce faire, on peut constater que pour tous k, n fixés, l'événement $\{C_n > k\}$ est égal à $\bigcup_{i=1}^n A_i$ où

$$A_i = \{i \notin \{X_1, \dots, X_k\}\}.$$

La formule d'inclusion-exclusion donne alors que

$$\begin{aligned} \mathbb{P}(C_n \leq k) &= \sum_{r=0}^n (-1)^r \sum_{1 \leq i_1 < \dots < i_r \leq n} \mathbb{P}(A_{i_1} \cap \dots \cap A_{i_r}) \\ &= \sum_{r=0}^n (-1)^r \binom{n}{r} \left(1 - \frac{r}{n}\right)^k. \end{aligned}$$

Prenons maintenant $k = k_n = \lfloor n \ln(n) + nx \rfloor$ avec $x \in \mathbb{R}$. Alors pour tout $r \geq 0$, on a lorsque $n \rightarrow \infty$

$$\left(1 - \frac{r}{n}\right)^k = \exp(k \ln(1 - r/n)) = n^{-r} e^{-rx} (1 - o(1)),$$

où le signe $-$ devant $o(1)$ est là pour insister sur le fait que cette quantité positive est majorée par $n^{-r} e^{-rx}$ pour tout n , et donc

$$\mathbb{P}(C_n \leq n \ln n + nx) = \sum_{r=0}^n \frac{(-1)^r}{r!} \frac{n!}{(n-r)!} n^{-r} e^{-rx} (1 - o(1)).$$

Comme $n!/(n-r)! = n^r (1 - o(1))$, on déduit aisément

$$\mathbb{P}\left(\frac{C_n - n \ln n}{n} \leq x\right) \xrightarrow{n \rightarrow \infty} \sum_{r \geq 0} \frac{(-1)^r}{r!} e^{-rx} = \exp(-e^{-x}).$$

Ainsi, la suite de variables aléatoires $(C_n - n \ln n)/n$, $n \geq 1$ converge en loi vers la *loi de Gumbel* dont la fonction de répartition est

$$\mathbb{P}(\mathcal{G} \leq x) = \exp(-e^{-x}).$$

Cette loi apparaît naturellement dans des problèmes d'étude de maxima de variables aléatoires indépendantes (valeurs extrêmes). Noter la dissymétrie marquée de cette fonction, puisque l'on a

$$\exp(-e^{-x}) = 1 - e^{-x} (1 + o(1))$$

lorsque $x \rightarrow \infty$, tandis que $\exp(-e^{-x})$ tend vers 0 « doublement exponentiellement » lorsque $x \rightarrow -\infty$. Autrement dit, s'il est plausible que la collection soit complète au temps $n \ln n + 3n$ (probabilité de l'ordre de $1 - e^{-3}$), il est en revanche franchement invraisemblable qu'elle le soit au temps $n \ln n - 3n$ (probabilité de l'ordre de $\exp(-e^3)$). On parle d'une « transition abrupte » (*cutoff* en anglais).

L'algorithme d'Aldous-Broder n'est pas le meilleur connu en termes de rapidité d'exécution : un autre algorithme également très simple, dit *algorithme de Wilson*, le surpasse de ce point de vue [Wil96]. Néanmoins, pour notre propos, c'est l'algorithme d'Aldous-Broder qui donne le point de vue le plus simple sur l'arbre couvrant uniforme.

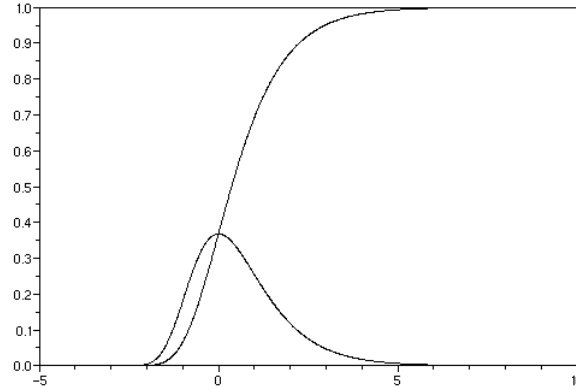


FIGURE 1. La fonction de répartition et la densité de la loi de Gumbel

2.3. Une autre interprétation

Comme nous l'avons dit, l'arbre couvrant uniforme du graphe complet que nous avons décrit ci-dessus n'est rien d'autre qu'un arbre uniforme parmi les n^{n-2} arbres possibles avec sommets étiquetés par $1, 2, \dots, n$. Il y a d'autres manières que l'algorithme d'Aldous-Broder d'envisager cet arbre d'un point de vue probabiliste, dont le recours aux arbres associés aux processus de branchement (Bienaymé-Galton-Watson). Rappelons que si μ est une mesure de probabilités sur \mathbb{N} , que l'on voit comme la loi de reproduction pour des individus d'une population asexuée, on peut lui associer un processus dans lequel chaque individu donne naissance indépendamment des autres à un nombre d'individus aléatoire de loi μ . À ce processus est à son tour associé l'arbre généalogique correspondant.

On renvoie au texte d'Igor Kortchemski (ce volume) pour la propriété suivante. Nous reviendrons brièvement sur cette interprétation au paragraphe 2.5.

Proposition 2.6. *Soit μ la loi de Poisson de paramètre $\lambda > 0$, et soit A l'arbre généalogique d'un processus de branchement de Bienaymé-Galton-Watson, de loi de reproduction μ , conditionné à avoir n sommets. On étiquette ses sommets uniformément au hasard. Le graphe sur $V = [n]$ obtenu est alors un arbre couvrant uniforme de K_n .*

2.4. Géométrie asymptotique : l'arbre continu aléatoire

Notons A_n une variable aléatoire qui est un arbre couvrant uniforme de K_n . On insiste sur la dépendance en n , ce dernier paramètre allant être envoyé vers l'infini. Que peut-on dire de la géométrie de A_n ?

2.4.1. Distance entre deux points marqués. L'algorithme d'Aldous-Broder permet de donner une réponse à cette question. Commençons par étudier la distance de graphe dans A_n entre deux sommets choisis uniformément au hasard (ou, ce qui revient au même, entre la racine 1 et un sommet choisi au hasard, pour des raisons claires de symétrie). Le corollaire 2.5 montre que cette distance de graphe a même loi que la première répétition $R_1 = R_1^{(n)}$ dans une suite X_1, X_2, \dots . Il y a un lien clair avec un problème très connu en probabilités, appelé le « paradoxe des anniversaires », qui consiste à observer que la probabilité qu'au moins deux individus dans un groupe de m personnes aient la même date d'anniversaire subit une transition abrupte de 0 à 1 lorsque m passe, disons, de 20 à 26.

Ici, le problème est le même : n est, disons, le nombre de dates d'anniversaire possibles (365 dans le véritable problème des anniversaires), et l'on peut voir $R_1^{(n)}$ comme le premier instant où, en observant les individus d'un groupe un à un, l'on s'aperçoit de l'existence d'une coïncidence d'anniversaires.

Proposition 2.7. *Lorsque $n \rightarrow \infty$, la variable aléatoire $R_1^{(n)}/\sqrt{n}$ converge en loi vers une variable de loi de Rayleigh, au sens où*

$$\mathbb{P}(R_1^{(n)} > r\sqrt{n}) \xrightarrow{n \rightarrow \infty} \exp(-r^2/2),$$

où il convient d'interpréter $\exp(-r^2/2)$ comme la queue de distribution de la loi à densité $re^{-r^2/2}dr\mathbf{1}_{\mathbb{R}_+}(r)$ (loi de Rayleigh).

Plus généralement, pour tout m on a la convergence en loi des m premiers temps de répétitions $n^{-1/2}(R_1^{(n)}, \dots, R_m^{(n)})$ vers un vecteur limite (L_1, \dots, L_m) , dont la loi est à densité :

$$\ell_1 \ell_2 \dots \ell_m \exp(-\ell_m^2/2) d\ell_1 d\ell_2 \dots d\ell_m \mathbf{1}_{\{0 < \ell_1 < \ell_2 < \dots < \ell_m\}}$$

Démonstration. Remarquons que si a est un entier,

$$\begin{aligned} \mathbb{P}(R_1^{(n)} > a) &= 1 \cdot (1 - 1/n) \cdot (1 - 2/n) \dots (1 - a/n) \\ &= \exp\left(\sum_{i=1}^a \ln(1 - i/n)\right). \end{aligned}$$

Si a est maintenant autorisé à dépendre de n , mais avec $a = O(\sqrt{n})$, on obtient

$$\mathbb{P}(R_1^{(n)} > a) = \exp\left(-n^{-1} \sum_{i=1}^a i + o(1)\right) = \exp(-a^2/2n + o(1)),$$

d'où le résultat. En fait, on peut être plus précis et constater que

$$\mathbb{P}(R_1^{(n)} = a) = \frac{a-1}{n} \mathbb{P}(R_1^{(n)} > a) = \frac{a}{n} \exp(-a^2/2n + o(1)).$$

Cet argument se généralise si l'on considère les répétitions suivantes : par exemple, pour $a < b$, avec $b = O(\sqrt{n})$,

$$\begin{aligned} \mathbb{P}(R_1^{(n)} = a, R_2^{(n)} = b) &= \mathbb{P}(R_2^{(n)} = b | R_1^{(n)} = a) \mathbb{P}(R_1^{(n)} = a) \\ &= (b/n) \cdot (1 - (a-1)/n) \cdot (1 - a/n) \dots (1 - (b-1)/n) \mathbb{P}(R_1^{(n)} = a) \\ &= (a/n)(b/n) \exp(-b^2/2n + o(1)). \end{aligned}$$

On déduit le résultat pour $m = 2$ par les caractérisations de la convergence en loi rappelés dans l'introduction, et ceci se généralise sans difficulté à m quelconque. Il convient néanmoins de montrer que l'expression de ℓ_1, \dots, ℓ_m donnée dans l'énoncé est bien une densité pour tout m , ce que nous laissons en exercice. \square

La loi jointe des variables aléatoires L_1, L_2, \dots apparaissant dans l'énoncé précédent peut se décrire très simplement en termes d'un processus de Poisson, c'est-à-dire une suite η_1, η_2, \dots de variables aléatoires telles que la suite $(\eta_i - \eta_{i-1}, i \geq 1)$ (avec la convention $\eta_0 = 0$) est formée de variables aléatoires i.i.d. de loi exponentielle de paramètre 1, c'est-à-dire $\mathbb{P}(\eta_1 \geq t) = \exp(-t)$. On a en effet que

$$(L_1, L_2, \dots) \stackrel{(\text{loi})}{=} (\sqrt{2\eta_1}, \sqrt{2\eta_2}, \dots).$$

Nous laissons cette preuve en exercice. On constate donc que les variables aléatoires L_1, L_2, \dots ont tendance à s'accumuler, au sens où les écarts $L_n - L_{n-1}$ convergent vers 0 presque sûrement, en contraste

avec les η_n . Pour s'en convaincre, on peut utiliser la loi des grands nombres, qui montre que $\eta_n \sim n$ presque sûrement, et écrire, avec la représentation ci-dessus,

$$L_n - L_{n-1} = \sqrt{2\eta_{n-1}} \left(\sqrt{1 + e_n/2\eta_{n-1}} - 1 \right) = O(\ln(n)/n^{1/2})$$

où $e_n = \eta_n - \eta_{n-1}$, et où l'on a utilisé le fait que

$$\limsup \frac{e_n}{\ln(n)} \leq 1, \quad \text{presque sûrement.}$$

Ce dernier fait est une conséquence aisée du lemme de Borel-Cantelli : on a en effet, comme $\mathbb{P}(e_i > x) = e^{-x}$,

$$\mathbb{P}(e_n > (1 + \varepsilon) \ln(n)) = n^{-(1+\varepsilon)},$$

qui est sommable, et donc presque sûrement $e_n \leq (1 + \varepsilon) \ln(n)$ pour tout n assez grand.

2.4.2. Construction de l'arbre continu. Comment interpréter ce résultat ? Rappelons que l'on peut voir A_n , dans la construction par l'algorithme d'Aldous-Broder, comme un processus de recollement de chaînes de longueurs $R_1^{(n)}, R_2^{(n)} - R_1^{(n)}, R_3^{(n)} - R_2^{(n)}, \dots$, la chaîne de longueur $R_{m+1}^{(n)} - R_m^{(n)}$ étant branchée sur un sommet uniforme sur l'arbre réduit $A_n[m]$. Nous savons que les différences de longueurs, renormalisées par \sqrt{n} , convergent lorsque $n \rightarrow \infty$ vers la suite $L_1, L_2 - L_1, L_3 - L_2, \dots$ de limite nulle que nous avons décrite ci-dessus. Ceci donne une intuition de ce que devrait être la structure limite de A_n , lorsque toutes les distances sont renormalisées par \sqrt{n} . La description suivante est appelée la *line-breaking construction* par Aldous [Ald91]. Elle consiste à brancher récursivement des segments réels de longueur $L_i - L_{i-1}, i \geq 1$, en recollant une des extrémités en un point uniformément choisi parmi la réunion des $i - 1$ premiers segments recollés.

Soit U_1, U_2, \dots une suite de variables aléatoires i.i.d. uniformes sur $[0, 1]$, indépendantes de L_1, L_2, \dots . Considérons la distance D_∞ sur \mathbb{R}_+ obtenue en « coupant » \mathbb{R}_+ aux points L_i , c'est-à-dire que $D_\infty(x, y) = |x - y|$ s'il existe i tel que $x, y \in [L_{i-1}, L_i[$, et $D_\infty(x, y) = \infty$ sinon, où l'on note $L_0 = 0$ par convention. On considère la plus petite relation d'équivalence \approx sur \mathbb{R}_+ telle que

$L_i \approx U_i L_i$ pour tout $i \geq 1$. Enfin, on considère le quotient métrique de \mathbb{R}_+ par \approx . Il s'agit de l'espace métrique obtenu en recollant les (paires de) points identifiés par \approx . Plus précisément, on note

$$D(x, y) = \inf \left\{ \sum_{i=1}^k D_\infty(x_i, y_i) \right\},$$

où l'infimum porte sur toutes les suites finies de points $x_1, y_1, \dots, x_k, y_k$ avec $x_1 = x$, $y_k = y$ et $y_i \approx x_{i+1}$ pour $1 \leq i \leq k-1$. Il n'est pas difficile de voir que cet infimum est en fait atteint, et que $D(x, y) = 0$ si et seulement si $x \approx y$. Le quotient métrique est l'espace quotient \mathbb{R}_+/\approx muni de la métrique induite par la (pseudo) distance D , encore notée D .

De façon intuitive, nous avons donc collé l'extrémité gauche de l'intervalle $[L_i, L_{i+1}[$ au point $U_i L_i$, qui appartient à la réunion des i premiers intervalles $[L_{j-1}, L_j[$, $1 \leq j \leq i$.

Définition 2.1. L'arbre continu aléatoire brownien est l'espace complété de $(\mathbb{R}_+/\approx, D)$. On le note \mathcal{T} .

Un autre nom communément utilisé est CRT, pour *Continuum Random Tree*. Pour énoncer quelques propriétés importantes de \mathcal{T} , donnons une définition.

Définition 2.2. Un espace métrique (M, d) est un \mathbb{R} -arbre si

- pour tout $x, y \in M$, il existe une application isométrique

$$\phi : [0, d(x, y)] \longrightarrow M \quad \text{avec } \phi(0) = x, \phi(d(x, y)) = y$$

(on dit que ϕ est une géodésique de x à y), et

- il n'existe pas d'application continue injective $\mathbb{S}^1 \rightarrow M$.

En particulier, si un espace (M, d) est un \mathbb{R} -arbre il est connexe (et même connexe par arcs) et sans cycle, au sens de la deuxième propriété caractéristique. On ajoute juste la propriété supplémentaire que la métrique est une métrique de longueur, c'est-à-dire que la distance entre deux points s'interprète comme la longueur d'un (plus court) chemin continu entre ces points. On peut penser à un \mathbb{R} -arbre comme à un recollement de copies isométriques de segments

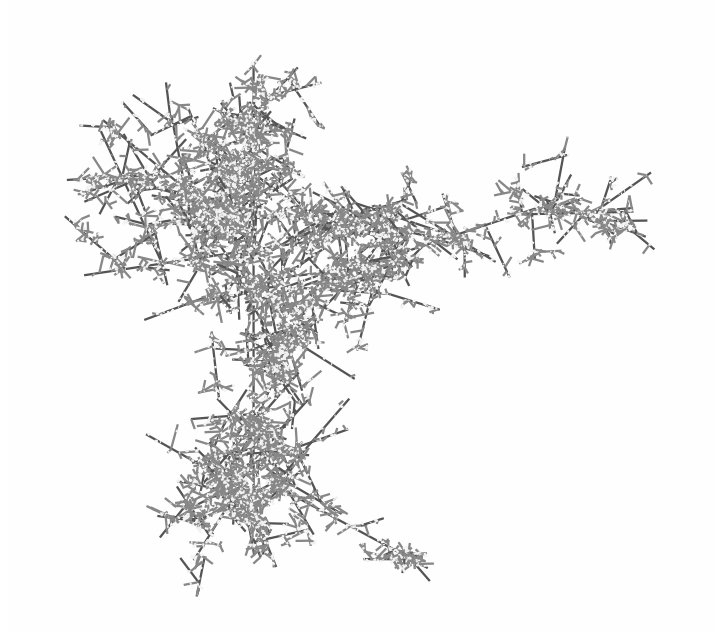


FIGURE 2. Quelques milliers d'itérations de la *line-breaking construction*, plongée dans \mathbb{R}^3 avec une orientation choisie uniformément au hasard pour chaque segment

réels, de façon à ne pas créer de cycle, comme le fait la *line-breaking construction*.

Théorème 2.8. *L'espace aléatoire \mathcal{T} est un \mathbb{R} -arbre presque sûrement, et sa dimension de Hausdorff est égale à 2 presque sûrement.*

Aldous [Ald91], montre la première partie de l'énoncé, ainsi qu'un résultat sur la dimension de Minkowski (le nombre de boules nécessaires pour recouvrir l'arbre), le résultat sur la dimension de Hausdorff apparaît dans Duquesne-Le Gall [DLG05].

Nous ne redonnons pas ici la définition de dimension de Hausdorff, mais insistons sur le côté un peu surprenant que peut avoir ce résultat à première vue : en effet, le CRT semble être un objet 1-dimensionnel, puisqu'il est obtenu à partir d'un recollement d'une quantité dénombrable de segments réels. Néanmoins, rappelons que \mathcal{T}

est en fait la *complétion* d'un tel recollement, et cette opération ajoute une quantité indénombrable de points (ces points sont nécessairement des *feuilles*, c'est-à-dire des points x tels que $\mathcal{T} \setminus \{x\}$ est connexe). Par analogie, on pourra penser aux ensembles de Cantor, que l'on obtient en retirant une collection dénombrable d'intervalles ouverts disjoints de $[0, 1]$: les extrémités gauches de ces intervalles forment un sous-ensemble dénombrable dense dans l'ensemble de Cantor correspondant, mais ce dernier est lui-même non dénombrable.

Enfin, nous avons indiqué dans quelle mesure on peut considérer que \mathcal{T} est un objet qu'il est naturel de considérer comme la limite d'échelle de l'arbre couvrant uniforme A_n . On peut donner un sens précis à cela : on dit que la suite d'espaces métriques compacts (M_n, d_n) converge vers une limite (M, d) au sens de Gromov-Hausdorff s'il existe un espace métrique (Z, δ) et des isométries $\phi_n : M_n \rightarrow Z$, $\phi : M \rightarrow Z$ telles que

$$\delta_H(\phi_n(M_n), \phi(M)) \xrightarrow{n \rightarrow \infty} 0,$$

où δ_H est la distance de Hausdorff associée à δ entre sous-ensembles fermés de Z , donnée par

$$\delta_H(A, B) = \max \left(\sup_{z \in A} d(z, B), \sup_{z \in B} d(z, A) \right).$$

On peut montrer [BBI01] que cette notion de convergence correspond effectivement à une topologie, et même une métrique, sur l'ensemble des espaces métriques compacts vus à isométrie près. De ce fait, cela a un sens de parler de convergence en loi de variables aléatoires à valeurs dans cet ensemble, muni de la distance de Gromov-Hausdorff.

Théorème 2.9 (Aldous [Ald91]). *On a que $(A_n, n^{-1/2}d_{A_n})$ converge vers \mathcal{T} en loi au sens de Gromov-Hausdorff.*

Aldous n'a pas énoncé ce résultat en termes de la distance de Gromov-Hausdorff, mais a en fait fourni une représentation (isométrique) de \mathcal{T} et des A_n renormalisés comme sous-ensembles compacts de $\ell^1(\mathbb{N})$, de sorte que la convergence a lieu en loi au sens de la distance de Hausdorff : cela implique immédiatement le résultat ci-dessus.

Les arguments que nous avons introduits impliquent assez facilement un résultat analogue, si nous remplaçons A_n par $A_n[m]$, et \mathcal{T} par le recollement $\mathcal{T}[m]$ des m premiers segments de la « line-breaking

construction », et ce pour tout m fixé. On voit que toute la difficulté pour passer d'un tel résultat à celui exposé ci-dessus est de montrer qu'on peut trouver m assez grand de sorte que $A_n[m]$ soit assez proche de A_n au sens de la distance de Hausdorff, et ce uniformément en n , au sens suivant : pour tout $\varepsilon > 0$,

$$\lim_{m \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbb{P}(\Delta(m, n) > \varepsilon \sqrt{n}) = 0$$

où $\Delta(m, n)$ est la distance maximale d'un sommet de A_n à un sommet de $A_n[m]$. Nous ne donnons pas ici la preuve de ce résultat technique, qui peut se démontrer de façon analogue aux *inégalités maximales* sur les marches aléatoires.

2.5. Perspectives sur le CRT

L'arbre brownien est un objet qui intervient un peu partout en probabilités, un peu au même titre que le mouvement brownien. Suivant l'approche originale d'Aldous, nous avons motivé sa définition par une approche algorithmique du problème de l'arbre couvrant uniforme du graphe complet, en insistant sur la *line-breaking construction* dont la description est peu technique — elle revient après tout à se donner les deux suites infinies (L_1, L_2, \dots) et (U_1, U_2, \dots) , dont la description est élémentaire.

Néanmoins, une définition plus moderne de l'arbre brownien, qui explique son nom, se fait en termes d'un objet bien plus élaboré, que l'on peut définir en termes du mouvement brownien $(B_t, t \geq 0)$ en dimension 1, et qui consiste à isoler une de ses excursions, c'est-à-dire la restriction de sa trajectoire à une des composantes connexes de l'ouvert $\{t \geq 0 : B_t > 0\}$. Nous ne rentrerons pas dans les détails de cette construction, issue des idées de Neveu–Pitman [NP89], Le Gall [LG93] et Aldous [Ald93], et qui est plus proche dans l'esprit de l'approche des arbres décrite dans le texte d'Igor Kortchemski (ce volume, [Kor16]). On pourra aussi consulter [LG05] pour aller plus loin.

La première raison qui explique le caractère naturel de l'arbre brownien est le fait qu'il peut se voir comme la limite universelle des processus de branchement, et non seulement la limite de l'arbre couvrant uniforme de K_n , qui, rappelons-le, peut se voir comme l'arbre

généalogique d'un processus de branchement dont la loi de reproduction est une loi de Poisson (et conditionné à avoir n individus).

Théorème 2.10 (Aldous [Ald93]). *Soit $\mu = (\mu(k), k \geq 0)$ une mesure de probabilités sur \mathbb{Z}_+ , telle que*

$$\mu(1) < 1, \quad \sum_{k \geq 0} k\mu(k) = 1, \quad \text{et} \quad \sigma^2 = \sum_{k \geq 0} k^2\mu(k) - 1 < \infty.$$

Soit A_n l'arbre généalogique d'un processus de branchement de Bienaymé-Galton-Watson de loi de reproduction μ , conditionné à avoir n sommets (on suppose que n est restreint aux entiers tels que le conditionnement a un sens.)

Alors $(A_n, (\sigma/\sqrt{n})d_{A_n})$ converge dans le même sens que le théorème 2.9 vers le CRT \mathcal{T} .

Par ailleurs, le CRT est présent, comme on l'a dit, dans l'étude asymptotique de très nombreux problèmes combinatoires. On pourra citer le problème de parking de Knuth [CL02], les arbres sur réseau [DS98], les processus de coalescence et de fragmentation [AP98, Ber06], les modèles d'attachement préférentiel, les partitions non croisées, les laminations discrètes du cercle (voir le texte d'Igor Kortchemski [Kor16]), et les cartes aléatoires [LGM12]. C'est peut-être dans ce dernier domaine, qui consiste à étudier de grands graphes aléatoires plongés dans le plan, que l'intervention du CRT est la plus surprenante, puisque ce dernier semble émerger d'une géométrie de type « champ moyen » (dans le graphe complet, tous les sommets sont voisins les uns des autres), alors que les cartes planaires sont des objets intrinsèquement 2-dimensionnels.

Intermezzo : le processus des forêts coalescentes de Pitman

Avant de changer de sujet, nous allons présenter une autre preuve de la formule de Cayley, due à Pitman, et qui jette un autre éclairage sur l'arbre couvrant uniforme de K_n . Une forêt sur $[n]$ est simplement un sous-graphe $([n], E)$ de K_n , qui est sans cycle mais plus nécessairement connexe. Une forêt est dite enracinée si chacune des composantes connexes de cette forêt, qui est évidemment un arbre, est enracinée. On notera que le nombre d'arbres (composantes connexes) d'une forêt $([n], E)$ est donné par $n - \text{Card}(E)$, puisque chaque arbre

a une arête de moins que de sommets, et tout sommet de $[n]$ est dans un arbre de la forêt.

On notera $\mathbf{F}_{n,k}$ l'ensemble des forêts enracinées sur $[n]$ à k arbres.

Théorème 2.11. *Pour tout $n \geq 1$ et $k \in \{1, 2, \dots, n\}$, on a*

$$\text{Card}(\mathbf{F}_{n,k}) = n^{n-k} \binom{n-1}{k-1}.$$

On récupère la formule de Cayley puisque cela donne

$$\text{Card}(\mathbf{F}_{n,1}) = n^{n-1},$$

et $\mathbf{F}_{n,1}$ est l'ensemble des arbres couvrants de K_n , enracinés. Comme chaque arbre sur $[n]$ a n choix possibles de racine, donnant autant d'arbres enracinés différents, on a que $\text{Card}(\mathbf{F}_{n,1}) = n \text{Card}(\mathbf{A}_n)$, et l'on récupère le théorème 2.1.

Nous allons montrer ce résultat par une construction combinatoire (et probabiliste) due à Pitman [Pit99]. Si $\mathbf{f} \in \mathbf{F}_{n,k}$, soit $\mathbf{t}_1, \dots, \mathbf{t}_k$ les composantes connexes de \mathbf{f} énumérées de façon arbitraire, les racines correspondantes étant désignées par r_1, \dots, r_k . Supposons que $k \geq 2$. Soit X une variable aléatoire uniforme dans $[n]$, et R une racine uniforme parmi celles des arbres de \mathbf{f} qui ne contiennent pas X : formellement, si $i \in [n]$, on a donc

$$\mathbb{P}(R = r_j | X = i) = \frac{1}{k-1}$$

pour tout $j \in \{1, 2, \dots, k\}$ tel que $i \notin V(\mathbf{t}_j)$. Notons \mathbf{f}' le graphe obtenu en ajoutant l'arête $\{X, R\}$. Clairement, \mathbf{f}' est une forêt enracinée à $k-1$ arbres, l'arbre obtenu par la fusion de celui contenant X et celui contenant R restant naturellement enraciné en la racine du premier (de sorte que R perd son rôle de sommet distingué dans \mathbf{f}'). On note $\mu(\mathbf{f}, \cdot)$ la loi de la forêt aléatoire \mathbf{f}' ainsi obtenue, qui est donc une mesure de probabilité sur les forêts enracinées.

Proposition 2.12. *Supposons que $F_1 = \mathbf{f}_1 = ([n], \emptyset)$ soit la forêt sur $[n]$ sans aucune arête (chaque arbre étant donc un singleton, naturellement enraciné). Soit F_2, \dots, F_n une suite de variables aléatoires*

obtenues successivement de F_1 par le procédé ci-dessus, ainsi, pour tout choix des forêts $\mathbf{f}_2, \dots, \mathbf{f}_n$, on a

$$\mathbb{P}(F_2 = \mathbf{f}_2, \dots, F_n = \mathbf{f}_n) = \mu(\mathbf{f}_1, \{\mathbf{f}_2\}) \cdots \mu(\mathbf{f}_{n-1}, \{\mathbf{f}_n\}),$$

cette probabilité étant nulle si $\mathbf{f}_k \notin \mathbf{F}_{n-k+1}$ pour un $k \in \{1, 2, \dots, n\}$.

Alors pour tout $k \in [n]$, F_k est uniforme parmi \mathbf{F}_{n-k+1} .

Démonstration. Le résultat est évidemment vrai pour $k = 1$. Supposons qu'il soit vrai pour l'indice $k - 1$. Soit \mathbf{f} une forêt enracinée à $n - k + 1$ arbres. Calculons la probabilité que $F_k = \mathbf{f}$. Pour que ceci soit vérifié, il faut que F_{k-1} soit égale à une forêt $\mathbf{f}(e)$ obtenue de \mathbf{f} en effaçant une arête e , et que le procédé décrit ci-dessus ait (ré)inséré cette arête. Notons $\mathbf{t}(e)$ l'arbre de \mathbf{f} contenant e et orientons naturellement $\vec{e} = (e_-, e_+)$ vers la racine $r(e)$ de $\mathbf{t}(e)$, $e_-, e_+ \in [n]$ étant donc l'origine et la cible de e . Alors $\mathbf{t}(e)$ privé de e consiste en deux composantes connexes $\mathbf{t}_+(e)$ et $\mathbf{t}_-(e)$, où $\mathbf{t}_+(e)$ contient e_+ et $r(e)$ et reste enraciné en $r(e)$, enraciné, tandis que $\mathbf{t}_-(e)$ est naturellement enraciné en la source e_- de \vec{e} . On a alors

$$\mathbb{P}(F_k = \mathbf{f}) = \sum_{e \in E(\mathbf{f})} \mathbb{P}(F_{k-1} = \mathbf{f}(e), X = e_+, R = e_-),$$

où X est uniforme dans $[n]$ et R est une racine uniforme parmi les $n - k + 1$ arbres de $\mathbf{f}(e)$ ne contenant pas X , comme ci-dessus. Ainsi,

$$\mathbb{P}(F_k = \mathbf{f}) = \sum_{e \in E(\mathbf{f})} \mathbb{P}(F_{k-1} = \mathbf{f}(e)) \cdot \frac{1}{n} \cdot \frac{1}{n - k + 1},$$

Or, par hypothèse de récurrence,

$$\mathbb{P}(F_{k-1} = \mathbf{f}(e)) = 1 / \text{Card}(\mathbf{F}_{n, n-k+2})$$

ne dépend pas de \mathbf{f} ni de e , et on a que $\text{Card}(E(\mathbf{f})) = k - 1$ D'où

$$(1) \quad \mathbb{P}(F_k = \mathbf{f}) = \frac{k - 1}{n(n - k + 1)\mathbf{F}_{n, n-k+2}}.$$

Cette quantité ne dépend que de n et de k , et pas de la valeur de $\mathbf{f} \in \mathbf{F}_{n, n-k+1}$, et l'on déduit le résultat. \square

La preuve du théorème 2.11 est alors immédiate, puisque (1) implique

$$\text{Card}(\mathbf{F}_{n,n-k+1}) = \frac{n(n-k+1)}{k-1} \text{Card}(\mathbf{F}_{n,n-k+2}),$$

et puisque l'on a $\text{Card}(\mathbf{F}_{n,n}) = 1$.

On propose enfin l'exercice suivant au lecteur intéressé. Supposons que $(p_1, p_2, \dots, p_n) \in \mathbb{R}_+^n$ soit un vecteur de probabilités sur $[n]$, c'est-à-dire que $p_1 + \dots + p_n = 1$. Effectuons la même construction que ci-dessus, mais supposons que dans la définition de μ , la variable aléatoire X soit choisie selon p , c'est-à-dire que $\mathbb{P}(X = i) = p_i$ pour tout $i \in [n]$. En revanche, on suppose toujours que la racine R est choisie uniformément parmi les racines des arbres ne contenant pas X . On définit comme avant la suite F_1, F_2, \dots, F_n , mais avec ce nouveau choix de la loi de X .

On rappelle que chaque arête e d'une forêt enracinée est naturellement orientée vers la racine de l'arbre la contenant, et l'on note $\vec{e} = (e_-, e_+)$. Montrer que l'on a, pour tout $n \geq 1$ et $k \in \{1, 2, \dots, n\}$, et tout $\mathbf{f} \in \mathbf{F}_{n,n-k+1}$

$$\mathbb{P}(F_k = \mathbf{f}) = \frac{1}{\binom{n-1}{k-1}} \prod_{e \in E(\mathbf{f})} p_{e_+},$$

ce que l'on peut encore noter, si $c_i(\mathbf{f}) = \text{Card}(\{e \in E(\mathbf{f}) : e_+ = i\})$ est le nombre « d'enfants » de $i \in [n]$ dans \mathbf{f} ,

$$\mathbb{P}(F_k = \mathbf{f}) = \frac{1}{\binom{n-1}{k-1}} \prod_{i \in [n]} p_i^{c_i(\mathbf{f})}.$$

En déduire que

$$\sum_{\mathbf{f} \in \mathbf{F}_{n,k}} \prod_{i \in [n]} p_i^{c_i(\mathbf{f})} = \binom{n-1}{k-1},$$

et plus généralement que

$$\sum_{\mathbf{f} \in \mathbf{F}_{n,k}} \prod_{i \in [n]} x_i^{c_i(\mathbf{f})} = \binom{n-1}{k-1} (x_1 + \dots + x_n)^{n-k},$$

l'égalité ayant lieu dans l'anneau des polynômes $\mathbb{Z}[x_1, \dots, x_n]$.

3. L'arbre couvrant minimal

On s'intéresse à présent à un modèle d'arbre couvrant aléatoire de K_n complètement différent, et qui correspond à ce que l'on appelle un « problème d'optimisation combinatoire ». Comme tout problème d'optimisation, il s'agit de minimiser une certaine fonction (« l'énergie ») définie sur un espace de configurations. Dans un problème d'optimisation « combinatoire », on insiste sur le fait que l'espace de configurations est fini, mais, typiquement, beaucoup trop grand pour espérer résoudre le problème par une fouille exhaustive des configurations. Beaucoup de ces problèmes peuvent se poser en termes d'un graphe pondéré $G = (V, E, w)$, c'est-à-dire un graphe muni d'une fonction $w : E \rightarrow (0, \infty)$ de poids sur les arêtes du graphe. Parmi ceux-ci :

- trouver l'arbre couvrant dont le poids $w(\mathbf{a}) = \sum_{e \in E(\mathbf{a})} w(e)$ est le plus petit possible : c'est le problème de l'arbre couvrant minimal.
- si (V, E) est le graphe complet à n sommets, trouver le cycle hamiltonien (visitant tous les sommets une fois et une seule), disons $C = (x_0, x_1, x_2, \dots, x_n = x_0)$ dont le poids

$$w(C) = \sum_{i=0}^{n-1} w(\{x_{i-1}, x_i\})$$

est le plus petit possible : c'est le problème du voyageur de commerce (*travelling salesman problem* en anglais).

Les deux problèmes présentés sont emblématiques et sont situés en quelque sorte aux deux extrémités du spectre en termes de difficulté : le problème du voyageur de commerce est *NP-hard*, c'est-à-dire « au moins aussi difficile que tout problème NP »⁽¹⁾. À l'inverse, le problème de l'arbre couvrant minimal est des plus simples, au sens où il peut se résoudre par des algorithmes *gloutons*, c'est-à-dire faisant

⁽¹⁾Un problème NP est un problème dont il est « facile » de vérifier si une solution donnée est effectivement une solution, c'est-à-dire qu'il existe un algorithme effectuant cette vérification en un temps polynomial en la longueur de la solution. En revanche, *trouver* une solution en un temps polynomial est beaucoup plus difficile, et, on le pense, en général impossible : c'est le fameux problème « P vs. NP ».

à chaque étape un choix d'optimum *local*, et fonctionnant de plus en temps polynomial. Néanmoins, le problème devient étonnamment riche lorsqu'on y ajoute une part d'aléa.

Reprenons le cas d'un graphe général fini connexe $G = (V, E)$. Soit donc $w(e), e \in E$ une famille de poids positifs indexés par les arêtes de G . Si \mathbf{a} est un arbre couvrant de G , on note

$$w(\mathbf{a}) = \sum_{e \in E(\mathbf{a})} w(e),$$

appelé poids de \mathbf{a} . Par la suite, on supposera *toujours* que les poids sont deux à deux distincts.

Définition 3.1. L'arbre couvrant minimal de G est l'arbre \mathbf{a} qui réalise le minimum des poids $w(\mathbf{a})$.

Il est clair que l'arbre couvrant minimal existe, puisqu'il n'y a qu'un nombre fini d'arbres couvrants. En revanche, on peut se demander s'il est bien défini de façon unique.

Proposition 3.1. *L'arbre couvrant d'un graphe pondéré dont les poids sont distincts deux à deux est unique.*

Démonstration. Supposons qu'il existe deux arbres couvrants minimaux distincts \mathbf{a}, \mathbf{a}' . Soit e l'arête de plus petit poids qui est dans l'un mais pas dans l'autre : on suppose sans perte de généralité que e est une arête de \mathbf{a} . Alors le graphe $\mathbf{a}' \cup \{e\}$ contient un cycle C . Par définition, toutes les arêtes de poids plus petit que $w(e)$ qui sont dans \mathbf{a}' sont aussi dans \mathbf{a} , et par conséquent, il existe une arête e' de C dont le poids est plus grand que celui de e (sinon toutes les arêtes de C seraient dans \mathbf{a} , mais ce dernier ne contient pas de cycle). Finalement, le graphe $(\mathbf{a}' \cup \{e\}) \setminus \{e'\}$ est un arbre couvrant, de poids $w(\mathbf{a}') - w(e') + w(e) < w(\mathbf{a}')$. Contradiction avec la minimalité de $w(\mathbf{a}')$. \square

3.1. L'algorithme de Kruskal, et applications

L'algorithme de Kruskal (qui, semble-t-il, remonte en fait à Boruvka) est un algorithme « glouton ». Nous l'énonçons dans le cas d'un graphe fini connexe pondéré général $G = (V, E, w)$. L'algorithme

consiste, partant du graphe vide, à ajouter une à une les arêtes du graphe par ordre croissant de poids, sauf lorsque l'ajout de l'arête en question crée un cycle. En pseudo-code, cela donne l'algorithme suivant.

Algorithme de Kruskal.

- (1) Numéroté les arêtes de $G = (V, E)$ en e_1, e_2, \dots de sorte que $w(e_1) < w(e_2) < \dots$.
- (2) Initialiser $S = Z = \emptyset$ et $I = 1$.
- (3) Tant que $S \neq V$ faire
 - Si $(V, Z \cup \{e_I\})$ est sans cycle, alors $S \leftarrow S \cup e_I, Z \leftarrow Z \cup \{e_I\}$
 - $I \leftarrow I + 1$.
- (4) Retourner Z .

Théorème 3.2. *L'algorithme de Kruskal produit l'unique arbre couvrant minimal de G pondéré par w .*

Démonstration. Soit $G = (V, E)$ le sous-graphe de K_n produit par l'algorithme. Par définition, G n'a aucun cycle, et donc définit une forêt couvrante du graphe complet, c'est-à-dire que chaque composante connexe est un arbre. Par ailleurs, soit V_1 l'ensemble des sommets de G connectés à 1, et supposons par l'absurde que $V_1 \neq [n]$. L'ensemble (non vide) des arêtes ayant une extrémité dans V_1 et l'autre dans $[n] \setminus V_1$ a un élément de plus petit poids, et clairement, cette arête aurait dû être ajoutée à l'étape de l'algorithme où elle est considérée. Donc G est un arbre couvrant de K_n .

Il reste à montrer que G est un arbre couvrant minimal. Pour cela, montrons par récurrence que pour tout $i \in \{0, 1, \dots, \text{Card}(V) - 1\}$, le graphe $G(i)$ construit au i -ème instant de l'algorithme où une arête est ajoutée peut être prolongé en un arbre couvrant minimal.

Pour $i = 0$ c'est évident en prenant n'importe quel arbre couvrant minimal. Si cela est vrai au rang $i < \text{Card}(V) - 1$, soit $G'(i)$ un arbre couvrant minimal contenant $G(i)$. Soit e la $i + 1$ -ème arête ajoutée dans l'algorithme de Kruskal. Si e est dans $G'(i)$ alors $G'(i)$ contient aussi $G(i + 1)$ et on peut poser $G'(i + 1) = G'(i)$. Sinon, l'ajout de e à $G'(i)$ crée un cycle C . Il existe forcément une arête $e' \neq e$ de ce cycle

qui n'est pas dans $G'(i)$: en effet, dans le cas contraire, l'apparition de e à l'étape $i+1$ aurait créé le cycle C , et e n'aurait pas été ajoutée à cette étape. Noter également que e' est forcément de poids plus grand que e : sinon, e' aurait été choisie par l'algorithme à une étape précédente, et e n'aurait pu être ajoutée sans créer de cycle. Alors $G'(i+1) = (G'(i) \setminus \{e'\}) \cup \{e\}$ est un arbre couvrant qui prolonge $G(i+1)$, et son poids est $w(G'(i)) + w(e) - w(e') < w(G'(i))$, ce qui est absurde par minimalité.

Cette preuve montre de plus que

$$G'(0) = G'(1) = \dots = G'(\text{Card}(V) - 1),$$

c'est-à-dire qu'il n'y a qu'un seul arbre couvrant minimal. \square

Nous en déduisons le critère suivant déterminant si une arête est dans l'arbre couvrant minimal, qui nous sera très utile par la suite.

Proposition 3.3. *Une arête $e = \{x, y\} \in E$ est dans $E(\mathbf{a}_m)$ si et seulement si pour tout chemin $x = x_1, x_2, \dots, x_k = y$, il existe une arête $e' = \{x_i, x_{i+1}\}$ telle que $w(e) \geq w(e')$.*

Démonstration. Supposons que $e = e_j$ soit la j -ème arête dans l'ordre croissant de poids. Cette arête est ajoutée à la j -ème étape de l'algorithme de Kruskal si et seulement si les extrémités x et y ne sont pas connectées à l'étape $j-1$. Si tout chemin de x à y passe par une arête de poids $> w(e)$, alors, comme aucune de ces arêtes ne peut être ajoutée avant l'étape j de l'algorithme, on voit que x et y ne peuvent être connectés à cette étape, et donc $e \in E(\mathbf{a}_m)$.

Réciproquement, s'il existe un chemin d'arêtes de $\{e_1, \dots, e_{j-1}\}$ connectant x et y , c'est-à-dire que x et y sont dans la même composante connexe C du graphe sur les arêtes $\{e_1, \dots, e_{j-1}\}$, alors l'algorithme de Kruskal à l'étape $j-1$ produit un arbre couvrant de C , et connecte donc x et y . Donc e n'est pas dans $E(\mathbf{a}_m)$. \square

Un autre algorithme, dû à Prim, consiste à ajouter l'arête de plus petit poids disponible, mais uniquement parmi celles incidentes à une « graine » fixée au départ. On fixe donc un sommet $v_0 \in V$.

Algorithme de Prim.

- (1) Initialiser $S = \{v_0\}$ et $Z = \emptyset$ et $I = 1$.

- (2) Tant que $S \neq V$ faire
- trouver l'arête $e = \{x, y\}$ de poids minimal avec $x \in S$ et $y \in V \setminus S$
 - $S \leftarrow S \cup \{y\}$ et $Z \leftarrow Z \cup \{e\}$
- (3) Retourner Z .

Nous laissons au lecteur le soin de démontrer que l'algorithme produit effectivement l'arbre couvrant minimal.

3.2. Le « théorème $\zeta(3)$ » de Frieze

Supposons maintenant que les poids $w(e), e \in E(K_n)$ soient des variables aléatoires i.i.d. de loi exponentielle de paramètre 1 :

$$\mathbb{P}(w(e) > r) = e^{-r} \quad \text{pour tout } r \geq 0.$$

La loi exponentielle étant à densité, il est effectivement vérifié presque sûrement que les poids sont deux à deux distincts, comme nous l'exigeons. Tacitement, on travaillera systématiquement en restriction à cet événement de probabilité 1.

On notera M_n l'arbre couvrant minimal associé à la suite des poids ci-dessus. On peut poser au moins deux questions naturelles au sujet de cet arbre :

- quelle est la géométrie asymptotique de M_n lorsque $n \rightarrow \infty$, au sens de la première partie (convergence de Gromov-Hausdorff)
- quel est le poids de M_n ?

On connaît essentiellement la réponse à ces deux questions, et nous allons évoquer la seconde avec quelque détail. Le résultat suivant est dû à Frieze en 1985 [Fri85].

Théorème 3.4. *On a $\mathbb{E}[w(M_n)] \rightarrow \zeta(3)$ lorsque $n \rightarrow \infty$, et de plus, pour tout $\varepsilon > 0$, on a que*

$$\mathbb{P}(|w(M_n) - \zeta(3)| > \varepsilon) \xrightarrow[n \rightarrow \infty]{} 0.$$

Il peut paraître surprenant de voir apparaître la quantité $\zeta(3)$ dans ce contexte. Mais la première surprise vient certainement du fait qu'il n'y a aucune renormalisation dans la convergence, alors que les poids sont des variables aléatoires exponentielles, et que l'arbre minimal

contient $n - 1$ arêtes (comme tout arbre couvrant de K_n). Néanmoins, le mystère est dissipé par le résultat suivant, qui sera utile par la suite.

Proposition 3.5. *Soit e_1, e_2, \dots, e_n des variables aléatoires indépendantes exponentielles de paramètre 1. Notons $e_{(1)} < e_{(2)} < \dots < e_{(n)}$ le réordonnement croissant de $\{e_1, \dots, e_n\}$. Alors $(e_{(1)}, \dots, e_{(n)})$ a même loi que la suite des sommes partielles renormalisées*

$$\sum_{i=1}^j \frac{e_i}{n-i+1}, \quad 1 \leq j \leq n$$

Démonstration. Soit $f : \mathbb{R}^n \rightarrow \mathbb{R}_+$ une fonction mesurable. Alors

$$\begin{aligned} \mathbb{E}[f(e_{(1)}, \dots, e_{(n)})] \\ = n! \int_{0 \leq x_1 \leq \dots \leq x_n} dx_1 \dots dx_n e^{-(x_1 + \dots + x_n)} f(x_1, \dots, x_n). \end{aligned}$$

On fait le changement de variables $y_i = x_i - x_{i-1}$, $1 \leq i \leq n$ (avec la convention $x_0 = 0$), ce qui donne

$$\begin{aligned} \mathbb{E}[f(e_{(1)}, \dots, e_{(n)})] \\ = n! \int_{\mathbb{R}_+^n} dy_1 \dots dy_n e^{-\sum_{i=1}^n (n-i+1)y_i} f(y_1, y_1 + y_2, \dots, y_1 + \dots + y_n), \end{aligned}$$

et l'on conclut en utilisant le fait que e_1/k a une loi à densité $ke^{-kx} \mathbf{1}_{\{x \geq 0\}}$. \square

On voit en particulier que tout sommet du graphe complet, qui est incident à $n - 1$ arêtes, est typiquement incident à au moins une arête de poids $O(1/n)$, et ce sont certainement ces très petites arêtes qui seront préférées par l'algorithme de Kruskal. Ceci explique l'absence de renormalisation dans le théorème de Frieze.

Nous allons donner l'idée d'une approche introduite par Aldous et Steele [AS92, AS04] pour calculer la limite de $\mathbb{E}[w(M_n)]$. Cette méthode, appelée la *méthode objective*, consiste à construire une limite (une sorte de version $n = \infty$) des objets étudiés. L'approche est longue à mettre en place, mais elle introduit en chemin des objets fondamentaux qui permettent d'obtenir une meilleure intuition du résultat et des objets avec lesquels on travaille. Nous insisterons davantage sur ces objets en eux-mêmes que sur les détails de la preuve.

Pour mettre en action la méthode objective, on est amené à étudier les propriétés *locales* de l'arbre couvrant du graphe complet, c'est-à-dire des propriétés que l'on peut décrire en étudiant seulement un voisinage d'une racine choisie au hasard.

Pourquoi une telle approche a-t-elle une chance de marcher ? Constatons que

$$\begin{aligned} \mathbb{E}[w(M_n)] &= \mathbb{E}\left[\sum_{e \in E(M_n)} w(e)\right] \\ &= \frac{1}{2} \mathbb{E}\left[\sum_{x \in [n]} \sum_{y \in [n]} w(\{x, y\}) \mathbf{1}_{\{\{x, y\} \in E(M_n)\}}\right] \\ &= \frac{1}{2n} \sum_{x \in [n]} \mathbb{E}\left[\sum_{y: \{x, y\} \in E(M_n)} nw(\{x, y\})\right] \\ &= \frac{1}{2} \mathbb{E}\left[\sum_{y: \{x_*, y\} \in E(M_n)} nw(\{x_*, y\})\right], \end{aligned}$$

où x_* est une variable aléatoire uniforme dans $[n]$, indépendante de $w(e)$, $e \in E(K_n)$. Nous avons donc exprimé l'espérance de la variable « globale » $w(M_n)$ comme celle de la variable « locale » (autour d'un sommet x_* choisi au hasard) $\sum_{y \sim x_*} nw(\{x_*, y\})$. L'idée est alors de trouver une forme de limite lorsque $n \rightarrow \infty$ de cette quantité. En fait, pour des raisons de symétrie, on peut remplacer x_* par 1, et nous sommes amenés à étudier la convergence en loi de

$$(2) \quad \frac{1}{2} \sum_{y: \{1, y\} \in E(M_n)} nw(\{1, y\}).$$

Remarquons que la seule convergence en loi d'une suite de variables aléatoires X_n positives ne permet pas de montrer la convergence de l'espérance : il faut pour cela montrer un résultat d'uniforme intégrabilité, c'est-à-dire que $\sup_n \mathbb{E}[X_n \mathbf{1}_{\{X_n > a\}}] \rightarrow 0$ lorsque $a \rightarrow \infty$. Nous admettrons que ceci est vérifié pour les variables définies en (2).

3.3. La convergence locale

La méthode objective utilise comme principal outil la notion de *convergence locale* d'une suite de graphe, que l'on peut exprimer

ainsi : on dit que la suite de graphes G_n à n sommets converge localement vers un graphe enraciné infini aléatoire (G, o) si, en choisissant o_n uniforme parmi les sommets de G , on a que pour tout r , la boule de rayon r centrée en o_n dans G_n converge en loi vers la boule de centre o et de rayon r dans G (pour la distance de graphe, et on voit $B_r(G, o)$ comme un graphe en y adjoignant toutes les arêtes de G reliant deux sommets à distance au plus r de o), c'est-à-dire que pour tout graphe enraciné (g, ρ) on a

$$\mathbb{P}(B_r(G_n, o_n) \approx (g, \rho)) \xrightarrow{n \rightarrow \infty} \mathbb{P}(B_r(G, o) \approx (g, \rho)),$$

où la notation \approx signifie que les graphes (enracinés) de part et d'autre sont isomorphes.

On peut se demander ce que signifie un « graphe infini aléatoire », puisque cela nécessite de préciser une notion de tribu : en fait, il est facile de voir que la convergence ci-dessus correspond à la convergence en loi pour une métrique sur les graphes connexes enracinés ayant une infinité dénombrable de sommets, par exemple

$$d_{\text{loc}}((G, o), (G', o')) = (1 + \sup\{r \geq 0 : B_r(G, o) \approx B_r(G', o')\} + 1)^{-1}.$$

On considère la tribu borélienne associée.

Si maintenant $G = (V, E, w)$ est pondéré, on peut modifier la définition précédente, en remplaçant les boules combinatoires de rayon r par les boules de rayon r dans la métrique induite par $w : d_w(x, y) = \inf\{\sum_{e \in \gamma: x \rightarrow y} w(e)\}$ où l'infimum est pris sur tous les chemins de x à y . On notera $B_r(G, o, w)$ la boule de rayon r pour cette distance, ce que l'on voit comme un graphe (non pondéré, mais que l'on peut munir naturellement de la restriction de w).

On dit que la suite de graphes pondérés enracinés (G_n, o_n, w_n) converge vers le graphe infini enraciné $(G = (V, E), o, w)$ si pour tout $r \notin \{d_w(o, v) : v \in V\}$, et pour tout n assez grand, il existe un isomorphisme de graphes $\phi_{n,r} : B_r(G, o) \rightarrow B_r(G_n, o_n)$, tel que $w_n(\phi_{n,r}(e))$ converge vers $w(e)$ pour tout e arête de $B_r(G, o, w)$. La raison pour laquelle on exclut le cas où $r = d_w(o, v)$ pour un $v \in V$ est que les boules de rayon s proche de r peuvent inclure ou exclure v selon que $s < r$ ou $s > r$: on ne peut donc pas exiger d'une suite de graphes approchant G que ses boules de rayon r soient toutes isomorphes.

Ici encore on peut munir l'ensemble des graphes enracinés pondérés infinis dénombrables d'une métrique correspondant à cette notion de convergence, par exemple

$$d_{\text{loc}}((G, o, w), (G', o', w')) = \int_0^\infty dr e^{-r} \left(\inf_{e \in E(B_r(G, o, w))} \sup_{\phi \in \text{Isom}_r} |w(e) - w'(\phi(e))| \wedge 1 \right),$$

où Isom_r est l'ensemble des isomorphismes de graphes entre

$$B_r(G, o, w) \quad \text{et} \quad B_r(G', o', w'),$$

ce dernier pouvant être vide.

Définition 3.2. On dit alors que la suite de graphes pondérés (éventuellement aléatoires) $(G_n, w_n), n \geq 1$, tel que G_n a n sommets, converge localement vers un graphe infini pondéré enraciné aléatoire (G, o, w) si, o_n étant un sommet choisi uniformément dans $V(G_n)$ conditionnellement à G_n , on a que (G_n, o_n, w_n) converge en loi vers (G, o, w) pour la distance d_{loc} .

La notion de convergence locale a été formalisée par Benjamini et Schramm [BS01], même si elle apparaît en germe dans des travaux antérieurs. Son importance dans l'étude de grands graphes aléatoires est cruciale.

3.4. Le graphe complet pondéré vu comme un arbre

Revenons à notre but de décrire un espace limite au graphe complet pondéré. L'objet limite se décrit de la façon suivante. Notons \mathcal{U} l'ensemble des mots d'entiers

$$\mathcal{U} = \bigcup_{n \geq 0} \mathbb{N}^n,$$

où $\mathbb{N}^0 = \{\emptyset\}$. On notera $\mathcal{U}^* = \mathcal{U} \setminus \{\emptyset\}$. On peut interpréter \mathcal{U} comme un arbre infini enraciné en \emptyset , et où le sommet correspondant au mot (u_1, u_2, \dots, u_n) , souvent noté $u_1 u_2 \dots u_n$, est vu comme le u_n -ième enfant du mot (u_1, \dots, u_{n-1}) . En tant que graphe, \mathcal{U} est donc muni des arêtes $\{u, uk\}$ avec $u \in \mathcal{U}$ et $k \geq 1$. De façon équivalente, l'ensemble des arêtes s'identifie canoniquement à \mathcal{U}^* .

Donnons-nous à présent une famille $(\xi_u, u \in \mathcal{U}^*)$ de variables aléatoires, de telle sorte que pour tout $u \in \mathcal{U}$, la famille $(\xi_{uk}, k \geq 1)$ ait même loi qu'un processus de Poisson (η_1, η_2, \dots) tel qu'introduit au

paragraphe 2.7. On demande de plus que ces familles de variables aléatoires soient indépendantes entre elles.

On interprète $(\xi_u, u \in \mathcal{U})$ comme des longueurs d'arêtes, au sens où ξ_{uk} est la longueur de l'arête $e = \{u, uk\}$, et l'on notera indifféremment $\xi(e)$ cette quantité. Le graphe pondéré enraciné $\mathfrak{T} = (\mathcal{U}, \emptyset, \xi)$ a été baptisé *Poisson weighted infinite tree* par Aldous, ce qu'on pourrait traduire par *arbre infini avec pondération poissonnienne*. L'acronyme, PWIT, est semble-t-il une facétie de l'auteur.

Théorème 3.6. *La limite locale du graphe complet (K_n, nw) pondéré par des variables exponentielles indépendantes de paramètre 1, renormalisées par n , est le PWIT.*

Une chose remarquable est que ce théorème reste valable pour des poids i.i.d. qui suivent une loi admettant une densité f continue dans un voisinage à droite de 0, avec $f(0) = 1$ (ou $f(0) > 0$, mais cela nécessite une renormalisation par une constante multiplicative dans l'énoncé).

Nous ne donnerons pas une preuve complète de ce théorème, mais nous contenterons de quelques remarques élémentaires. Tout d'abord, pour des raisons évidentes de symétrie, le choix de la racine o_n n'importe pas, et peut être aussi bien choisi comme étant le sommet 1 ; ce que nous supposons.

Le PWIT est bien un graphe localement fini. Nous utiliserons ici et par la suite le fait suivant.

Lemme 3.7. *Si e_1, \dots, e_k sont des variables aléatoires exponentielles de paramètre 1 indépendantes, alors leur somme $e_1 + \dots + e_k$ a pour loi la loi Gamma(k) dont la densité par rapport à la mesure de Lebesgue*

$$\frac{x^{k-1} e^{-x}}{k!} dx \mathbf{1}_{\{x \geq 0\}}.$$

La preuve est une application simple de la formule de changement de variables, et est omise.

Soit $u = u_1 u_2 \dots u_n \in \mathcal{U}$ un mot. La loi de $d_\xi(u, \emptyset)$ celle de la somme des variables ξ_v où v décrit l'ensemble des préfixes de u , et ces variables aléatoires sont indépendantes, respectivement de lois

$\text{Gamma}(u_i)$, $1 \leq i \leq n$. Leur somme est donc de loi $\text{Gamma}(|u|_1)$ par une application du lemme 3.7. Or si X suit une loi $\text{Gamma}(k)$ on a

$$P(X < r) = \frac{1}{k!} \int_0^r x^{k-1} e^{-x} dx \leq r^k/k!.$$

Par conséquent, l'espérance du nombre de sommets à d_ξ -distance inférieure à r de l'origine (à l'exception de celle-ci) est

$$\begin{aligned} \mathbb{E} \left[\sum_{u \in \mathcal{U}^*} \mathbf{1}_{\{d_\xi(u, \emptyset) \leq r\}} \right] &\leq \sum_{u \in \mathcal{U}^*} \frac{r^{|u|_1}}{|u|_1!} = \sum_{k \geq 1} \sum_{l \geq 0} \frac{r^{k+l}}{(k+l)!} \\ &\leq \sum_{k, l \geq 0} \frac{r^{k+l}}{k!l!} \leq e^{2r}, \end{aligned}$$

puisque $(k+l)! \geq k!l!$. D'où le résultat.

La première génération du PWIT. Nous avons donné précédemment la loi du réarrangement croissant $(e_{(1)}^n, \dots, e_{(n-1)}^n)$ de $n-1$ variables exponentielles indépendantes : on peut les représenter sous la forme

$$e_{(i)}^n = \sum_{j=1}^i \frac{e'_j}{n-j}, \quad 1 \leq i \leq n-1,$$

où e'_i , $i \geq 1$ sont également des variables exponentielles indépendantes. Par conséquent, on voit immédiatement que dans cette représentation,

$$n(e_{(i)}^n, i \geq 1) \xrightarrow{n \rightarrow \infty} \left(\sum_{j=1}^i e'_j, i \geq 1 \right),$$

(presque sûrement !) et on reconnaît à droite un processus de Poisson. Si l'on interprète les variables $(e_{(i)}^n, 1 \leq i \leq n)$ comme les distances aux voisins de la racine dans le graphe complet, on voit donc qu'elles correspondent à la limite aux distances de la racine à ses voisins dans le PWIT. On notera que cette limite décrit en réalité les voisins de la racine qui en sont très proches : la plupart des voisins se trouvent à une distance renormalisée de l'ordre de n , et « disparaissent » à la limite.

La preuve du théorème 3.6 consiste essentiellement à réitérer l'argument ci-dessus génération par génération, en justifiant le fait que les distances des B plus proches voisins de la racine de K_n à leurs B

plus proches voisins distincts de la racine peuvent être décrits à la limite par les B premiers atomes d'un processus de Poisson indépendants, et aussi que tous les sommets ainsi visités sont tous distincts. Ici, B est un grand entier fixé lorsque $n \rightarrow \infty$. Nous omettons les détails techniques de la preuve.

3.5. L'arbre minimal vu comme une forêt couvrante du PWIT

Rappelons que nous voulons décrire non seulement une limite locale du graphe complet, mais aussi de l'arbre couvrant minimal M_n . On peut se demander à quel objet cela peut correspondre sur le PWIT, puisque ce dernier est *déjà un arbre!* Mais un moment de réflexion montre que l'on ne saurait voir toutes les arêtes du PWIT comme correspondant à la limite de l'arbre couvrant minimal, pour la raison suivante. Rappelons-nous la proposition 3.3 disant qu'une arête $e = \{x, y\}$ est dans l'arbre couvrant minimal si et seulement si tout chemin de x à y emprunte au moins une arête e' de poids $w(e') > w(e)$. Ou autrement dit, si les extrémités de e deviennent déconnectées si l'on enlève toutes les arêtes de poids $\geq w(e)$.

Cette remarque permet de définir l'analogie de l'arbre couvrant minimal dans le PWIT : on voudrait dire que l'arête $e = \{u, uk\}$ est dans la *forêt couvrante minimale* $F(\mathfrak{T})$ si u et uk ne sont pas connectés dans le sous-graphe du PWIT obtenu en retirant toutes les arêtes de poids $\geq \xi(uk)$. Évidemment, cela est encore vide de sens, puisque le PWIT est un arbre, et donc deux sommets ne peuvent être connectés que par un seul chemin. Néanmoins, on peut rattraper cela en disant que les sommets u et uk ne sont pas connectés (hors de l'arête e) si l'une au moins des deux composantes connexes de ces deux sommets, formée des arêtes du PWIT de poids $< \xi(e)$ est finie. À l'inverse, on dira que l'arête e n'est pas dans $F(\mathfrak{T})$ s'il existe deux chemins infinis issus de u et uk et n'empruntant que des arêtes de poids $< \xi(e)$, l'idée étant que dans ce cas, u et uk sont connectés « à l'infini » par un chemin de tels arêtes. Nous allons expliquer pourquoi cette idée a une chance de fonctionner, mais tout d'abord nous donnons le résultat.

Définition 3.3. Soit $\mathfrak{T} = (\mathcal{U}, \emptyset, \xi)$ le PWIT. La forêt minimale couvrante $F(\mathfrak{T})$ de \mathfrak{T} est le sous-graphe de \mathfrak{T} défini par les arêtes $e = \{u, v\} \in E(\mathfrak{T})$ telles que l'une au moins des deux composantes connexes du sous graphe $(V(\mathfrak{T}), \{e' \in E(\mathfrak{T}) : \xi(e') < \xi(e)\})$ contenant u et v est finie.

On a alors le résultat suivant.

Théorème 3.8. *On a la convergence locale suivante*

$$(K_n, M_n, nw) \xrightarrow[n \rightarrow \infty]{} (\mathfrak{T}, F(\mathfrak{T}), \xi)$$

Il faut expliquer ce qu'on entend dans le théorème précédent par la convergence jointe de K_n et son arbre couvrant minimal M_n . Ce dernier est un sous-graphe de K_n , il est donc immédiat d'étendre la notion de convergence locale à (K_n, M_n, w) , si l'on voit par exemple M_n comme une fonction de marques $m : E(K_n) \rightarrow \{0, 1\}$, les pondérations w servant toujours à définir les distances : dans la définition ci-dessus de la distance locale, il faut se restreindre aux isomorphismes $\phi : B_r(G, o, w) \rightarrow B_r(G', o', w')$ qui préservent les marques : $m(e) = m'(\phi(e))$ pour tout $e \in E(B_r(G, o, w))$.

Le théorème 3.8 sera admis. On peut penser qu'il est un peu miraculeux que deux sommets voisins incidents à l'arête e du PWIT, parce qu'ils sont reliés à l'infini par des chemins n'empruntant pas e faits d'arêtes de poids $< \xi(e)$, peuvent effectivement être interprétés comme étant reliés ensemble. La raison est que, s'il existe un tel chemin, il en existe en fait énormément, et ces chemins envahissent littéralement le graphe. On a en effet le résultat suivant.

Lemme 3.9. *La composante connexe \mathfrak{T}_λ du sous-graphe*

$$(V(\mathfrak{T}), \{e \in E(\mathfrak{T}) : \xi(e) < \lambda\})$$

contenant \emptyset est l'arbre généalogique d'un processus de branchement de loi de reproduction Poisson de moyenne λ .

Démonstration. La propriété fondamentale du processus de Poisson, encore noté (η_1, η_2, \dots) est que $N_\lambda = \sup\{i \geq 0 : \eta_i < \lambda\}$ (avec

la convention $\eta_0 = 0$) est une variable aléatoire de Poisson de paramètre λ . En effet, pour $n \geq 0$, on a

$$\begin{aligned} \mathbb{P}(N_\lambda = n) &= \mathbb{P}(\eta_n \leq \lambda < \eta_{n+1}) \\ &= \int_{x_1 + \dots + x_n \leq \lambda < x_1 + \dots + x_n + x_{n+1}} dx_1 \dots dx_n dx_{n+1} e^{-(x_1 + \dots + x_n + x_{n+1})} \\ &= \int_{y > \lambda} dy e^{-y} \int_{y_1 < \dots < y_n \leq \lambda} dy_1 \dots dy_n \\ &= \frac{\lambda^n e^{-\lambda}}{n!}, \end{aligned}$$

comme voulu. Par construction du PWIT, on voit donc que le sous graphe de \mathfrak{T} formé des arêtes de poids $< \lambda$ est constitué des N_λ^u premiers voisins de chaque sommet $u \in V(\mathfrak{T})$, où N_λ^u , $u \in V(\mathfrak{T})$ est une famille i.i.d. de loi de Poisson de paramètre λ . Il est immédiat de vérifier la composante connexe de \emptyset a la loi voulue. \square

Les résultats standards sur le processus de branchement (rappelons que $\mathbb{E}[N_\lambda] = \lambda$) montrent que \mathfrak{T}_λ est fini presque sûrement si $\lambda \in [0, 1]$, et que $\mathbb{P}(|\mathfrak{T}_\lambda| = \infty)$ est l'unique racine $q = q(\lambda)$ dans $]0, 1[$ de l'équation

$$e^{-\lambda q} = 1 - q$$

(voir par exemple le texte d'Igor Kortchemski, ce volume). En fait, il est également connu que sur l'événement où \mathfrak{T}_λ est infini, la génération n contient de l'ordre de $W\lambda^n$ individus, où $W > 0$ est une variable aléatoire. On voit donc qu'il existe une quantité de chemins vers l'infini qui croît de façon exponentielle avec la génération : en fait, si l'approximation de K_n par le PWIT est bonne jusqu'à une hauteur de l'ordre de $c \ln(n)$ pour une constante $c > 0$, on voit qu'une proportion positive des sommets de K_n est reliée aux sommets correspondant à u et uk par des arêtes de poids $< \xi(e)$. Il devient alors plausible que l'on puisse en fait relier ces sommets dans K_n par un tel chemin.

Du théorème 3.8 on en déduit la convergence en loi

$$\frac{1}{2} \sum_{i: \{1, i\} \in E(M_n)} nw(\{1, i\}) \xrightarrow{n \rightarrow \infty} \frac{1}{2} \sum_{i: \{\emptyset, i\} \in F(\mathfrak{T})} \xi(i),$$

où la quantité à gauche est celle qui apparaît en (2). Comme nous avons admis que nous avons aussi la convergence des espérances,

il nous reste à calculer

$$\mathbb{E} \left[\sum_{i: \{\emptyset, i\} \in F(\mathfrak{T})} \xi(i) \right] = \mathbb{E} \left[\sum_{i \geq 1} \xi(i) \mathbf{1}_{\{|\mathfrak{T}_{\xi(i)}| < \infty \text{ ou } |\mathfrak{T}_{\xi(i)}^i| < \infty\}} \right],$$

où $\mathfrak{T}_{\xi(i)}$ est \mathfrak{T}_λ pour la valeur $\lambda = \xi(i)$, et $\mathfrak{T}_{\xi(i)}^i$ est la composante connexe contenant le sommet $i \in \mathcal{U}$ de $\{e \in E(\mathfrak{T}) : \xi(e) < \xi(i)\}$. Notons que les événements $\{|\mathfrak{T}_{\xi(i)}^i| = \infty\}$ et $\{|\mathfrak{T}_{\xi(i)}| = \infty\}$ sont indépendants conditionnellement à $\{\xi(u) : u \notin i\mathcal{U}^*\}$ où u décrit tous les sommets du PWIT, à l'exception des descendants stricts de i , car ils font intervenir des poids d'arêtes sur des générations disjointes du PWIT (à l'exception du poids $\xi(i)$ de l'arête $\{\emptyset, i\}$), et de plus la probabilité conditionnelle du premier est $\mathbb{P}(|\mathfrak{T}_{\xi(i)}^i| = \infty \mid \xi(u), u \notin i\mathcal{U}^*) = q(\xi(i))$, où $q(\lambda)$, défini au-dessus, est la probabilité de non-extinction de \mathfrak{T}_λ . On a donc

$$(3) \quad \mathbb{E} \left[\sum_{i: \{\emptyset, i\} \in F(\mathfrak{T})} \xi(i) \right] = \mathbb{E} \left[\sum_{i \geq 1} \xi(i) \left(1 - q(\xi(i)) \mathbf{1}_{\{|\mathfrak{T}_{\xi(i)}| = \infty\}} \right) \right].$$

Nous utilisons une dernière propriété des processus de Poisson : qui est que

$$\mathbb{E} \left[\sum_{i \geq 1} H(\eta_i, (\eta_j, j \geq 1, j \neq i)) \right] = \int_0^\infty d\lambda \mathbb{E}[H(\lambda, (\eta_j, j \geq 1))],$$

où H est une fonction mesurable positive (pouvant aussi être fonction de variables aléatoires indépendantes de η_1, η_2, \dots). On laisse au lecteur le soin de montrer cette propriété dans le cas où $H(\lambda, (t_j, j \geq 1))$ est de la forme $f(\lambda)g(t_1, \dots, t_k)$, le résultat s'en déduisant par un argument de classe monotone.

L'application de cette formule à (3) donne

$$\begin{aligned} \frac{1}{2} \mathbb{E} \left[\sum_{i: \{\emptyset, i\} \in F(\mathfrak{T})} \xi(i) \right] &= \int_0^\infty \lambda (1 - \mathbb{P}(|\mathfrak{T}_\lambda| = \infty)) q(\lambda) d\lambda \\ &= \frac{1}{2} \int_0^\infty \lambda (1 - q(\lambda))^2 d\lambda \\ &= \frac{1}{2} \int_1^\infty \lambda^2 q(\lambda) q'(\lambda) d\lambda \end{aligned}$$

où l'on a utilisé une intégration par parties, puis le fait que q est constante égale à 1 sur $[0, 1]$. Comme on a

$$\exp(-\lambda q(\lambda)) = (1 - q(\lambda)),$$

soit $\lambda = -\ln(1 - q(\lambda))/q(\lambda)$, on peut faire le changement de variables $q = q(\lambda)$, ce qui donne, avec un dernier changement de variables $u = -\log(1 - q)$,

$$\begin{aligned} \frac{1}{2} \int_1^\infty \lambda^2 q(\lambda) q'(\lambda) d\lambda &= \int_0^1 \frac{\log^2(1 - q)}{2q} dq = \int_0^1 \frac{u^2}{2} \frac{e^{-u}}{1 - e^{-u}} du \\ &= \int_0^1 \frac{u^2}{2} \sum_{k \geq 1} e^{-ku} du = \sum_{k \geq 1} \frac{1}{k^3} = \zeta(3). \end{aligned}$$

C'est ce qu'on voulait !

3.6. Géométrie asymptotique

On peut légitimement se poser la question de la convergence au sens de Gromov-Hausdorff de l'arbre minimal M_n , après renormalisation idoine. C'est l'objet d'un résultat récent [ABBGM17].

Théorème 3.10. *La suite $(M_n, n^{-1/3} d_{M_n})$ converge en loi au sens de Gromov-Hausdorff vers une limite \mathcal{M} , qui est un \mathbb{R} -arbre, dont la dimension de Minkowski égale 3 presque sûrement.*

On voit en particulier que la géométrie asymptotique de M_n est très différente de celle du CRT, dont la dimension (de Hausdorff, ou de Minkowski) est 2. En fait, on peut décrire la limite \mathcal{M} en termes d'une opération un peu sophistiquée sur le CRT, mais nous n'entrerons pas dans ces considérations.

Ce résultat fait suite aux importants travaux [ABBG12, ABBG10] dont l'objet est de décrire la limite d'échelle du *graphe d'Erdős-Rényi* [ER60] au paramètre critique, qui consiste à garder chaque arête du graphe complet avec probabilité p , et à enlever les autres. Les composantes connexes de ce graphe, lorsque p est de l'ordre de $1/n$, sont étroitement liées au CRT, tandis que l'émergence de ces composantes alors que p augmente est lié au déroulement de l'algorithme de Kruskal.

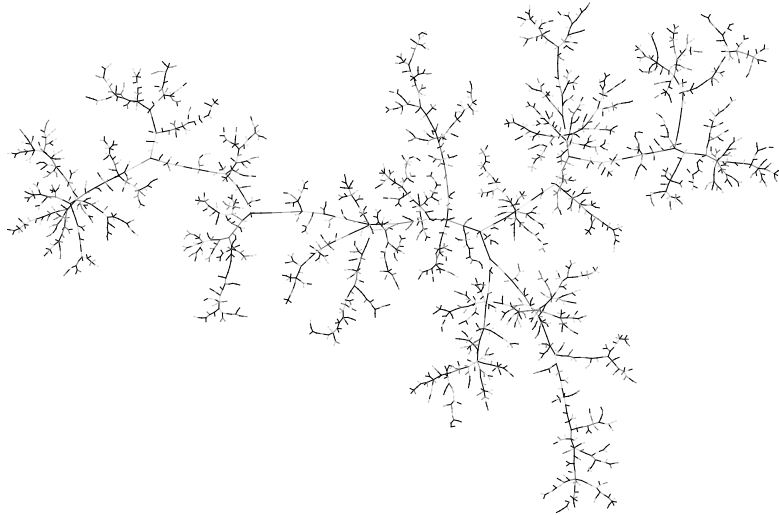


FIGURE 3. Une simulation de l'arbre couvrant minimal M_n avec $n = 3000$, par Louigi Addario-Berry

Le modèle du graphe d'Erdős-Rényi est un des modèles aléatoires les plus étudiés et utilisés dans la littérature scientifique, et il mériterait un ou plusieurs cours à lui tout seul. On consultera par exemple le très beau livre de van der Hofstad, *Random graphs and complex networks*, encore en préparation mais disponible en ligne, ou [ABBG10] pour une mise en perspective (plus ardue) de certaines des approches présentées dans ces notes.

Références

- [ABBG10] L. ADDARIO-BERRY, N. BROUTIN & C. GOLDSCHMIDT – « Critical random graphs : limiting constructions and distributional properties », *Electron. J. Probab.* **15** (2010), p. 741–775, art. no. 25.
- [ABBG12] ———, « The continuum limit of critical random graphs », *Probab. Theory Relat. Fields* **152** (2012), no. 3-4, p. 367–406.
- [ABBG17] L. ADDARIO-BERRY, N. BROUTIN, C. GOLDSCHMIDT & G. MIERMONT – « The scaling limit of the minimum spanning tree of the complete graph », *Ann. Probab.* **45** (2017), no. 5, p. 3075–3144.
- [Ald90] D. ALDOUS – « The random walk construction of uniform spanning trees and uniform labelled trees », *SIAM J. Discrete Math.* **3** (1990), no. 4, p. 450–465.
- [Ald91] ———, « The continuum random tree. I », *Ann. Probab.* **19** (1991), no. 1, p. 1–28.

- [Ald92] ———, « Asymptotics in the random assignment problem », *Probab. Theory Relat. Fields* **93** (1992), no. 4, p. 507–534.
- [Ald93] ———, « The continuum random tree. III », *Ann. Probab.* **21** (1993), no. 1, p. 248–289.
- [AP98] D. ALDOUS & J. PITMAN – « The standard additive coalescent », *Ann. Probab.* **26** (1998), no. 4, p. 1703–1726.
- [AS92] D. ALDOUS & J. M. STEELE – « Asymptotics for Euclidean minimal spanning trees on random points », *Probab. Theory Relat. Fields* **92** (1992), no. 2, p. 247–258.
- [AS04] ———, « The objective method : probabilistic combinatorial optimization and local weak convergence », in *Probability on discrete structures*, Encyclopaedia Math. Sci., vol. 110, Springer, Berlin, 2004, p. 1–72.
- [BS01] I. BENJAMINI & O. SCHRAMM – « Recurrence of distributional limits of finite planar graphs », *Electron. J. Probab.* **6** (2001), article no. 23.
- [Ber06] J. BERTOIN – *Random fragmentation and coagulation processes*, Cambridge Studies in Advanced Math., vol. 102, Cambridge University Press, Cambridge, 2006.
- [BBI01] D. BURAGO, Y. BURAGO & S. IVANOV – *A course in metric geometry*, Graduate Studies in Math., vol. 33, American Mathematical Society, Providence, RI, 2001.
- [CP00] M. CAMARRI & J. PITMAN – « Limit distributions and random trees derived from the birthday problem with unequal probabilities », *Electron. J. Probab.* **5** (2000), article no. 2.
- [CL02] P. CHASSAING & G. LOUCHARD – « Phase transition for parking blocks, Brownian excursion and coalescence », *Random Structures Algorithms* **21** (2002), no. 1, p. 76–119.
- [DS98] E. DERBEZ & G. SLADE – « The scaling limit of lattice trees in high dimensions », *Comm. Math. Phys.* **193** (1998), no. 1, p. 69–104.
- [DLG05] T. DUQUESNE & J.-F. LE GALL – « Probabilistic and fractal aspects of Lévy trees », *Probab. Theory Relat. Fields* **131** (2005), no. 4, p. 553–603.
- [ER60] P. ERDŐS & A. RÉNYI – « On the evolution of random graphs », *Magyar Tud. Akad. Mat. Kutató Int. Közl.* **5** (1960), p. 17–61.
- [Fri85] A. M. FRIEZE – « On the value of a random minimum spanning tree problem », *Discrete Appl. Math.* **10** (1985), no. 1, p. 47–56.
- [Kor16] I. KORTCHEMSKI – « Arbres et marches aléatoires », in *Arbres et marches aléatoires*, Journées X-UPS, Les Éditions de l'École polytechnique, Palaiseau, 2016, ce volume.
- [LG93] J.-F. LE GALL – « The uniform random tree in a Brownian excursion », *Probab. Theory Relat. Fields* **96** (1993), no. 3, p. 369–383.
- [LG05] ———, « Random trees and applications », *Probab. Surv.* **2** (2005), p. 245–311.
- [LGM12] J.-F. LE GALL & G. MIERMONT – « Scaling limits of random trees and planar maps », in *Probability and statistical physics in two and more dimensions*, Clay Math. Proc., vol. 15, American Mathematical Society, Providence, RI, 2012, p. 155–211.
- [NP89] J. NEVEU & J. PITMAN – « The branching process in a Brownian excursion », in *Séminaire de Probabilités, XXIII*, Lect. Notes in Math., vol. 1372, Springer, Berlin, 1989, p. 248–257.
- [Pit99] J. PITMAN – « Coalescent random forests », *J. Combin. Theory Ser. A* **85** (1999), no. 2, p. 165–193.

- [Wil96] D. B. WILSON – « Generating random spanning trees more quickly than the cover time », in *Proceedings of the Twenty-eighth Annual ACM Symposium on the Theory of Computing (Philadelphia, PA, 1996)*, ACM, New York, 1996, p. 296–303.

Grégory Miermont, Université de Lyon, ÉNS de Lyon & Institut Universitaire de France

E-mail : gregory.miermont@ens-lyon.fr

Url : <http://perso.ens-lyon.fr/gregory.miermont/>