

# Fonctions $\zeta$ et topologie stochastique

Daniel Perez <sup>1,2,3</sup>

<sup>1</sup>École normale supérieure

<sup>2</sup>Université Paris-Saclay

<sup>3</sup>DataShape, INRIA

29 septembre 2021

- 1 Ensembles de surniveau et code-barres
- 2 Quantités importantes
- 3 Processus stochastiques
- 4 Conclusion

## Notation

Notons  $X$  un espace topologique connexe, localement connexe par arcs et compact et  $f : X \rightarrow \mathbb{R}$  une fonction continue.

On peut associer à  $f$  un arbre  $T_f$  de la manière suivante

## Théorème (Arbre de surniveau, P.)

La fonction  $d_f : X \times X \rightarrow \mathbb{R}$  définie par

$$d_f(x, y) := f(x) + f(y) - 2 \sup_{\gamma: x \mapsto y} \inf_{t \in [0, 1]} f \circ \gamma \quad (1.1)$$

est une pseudo-distance sur  $X$  et satisfait

$$d_f(x, y) = 0 \iff x, y \in \{f = r\} \text{ et même comp. connexe de } \{f \geq r\}. \quad (1.2)$$

De plus,  $T_f := X / \{d_f = 0\}$  muni de la topologie quotient est un arbre réel (enraciné).

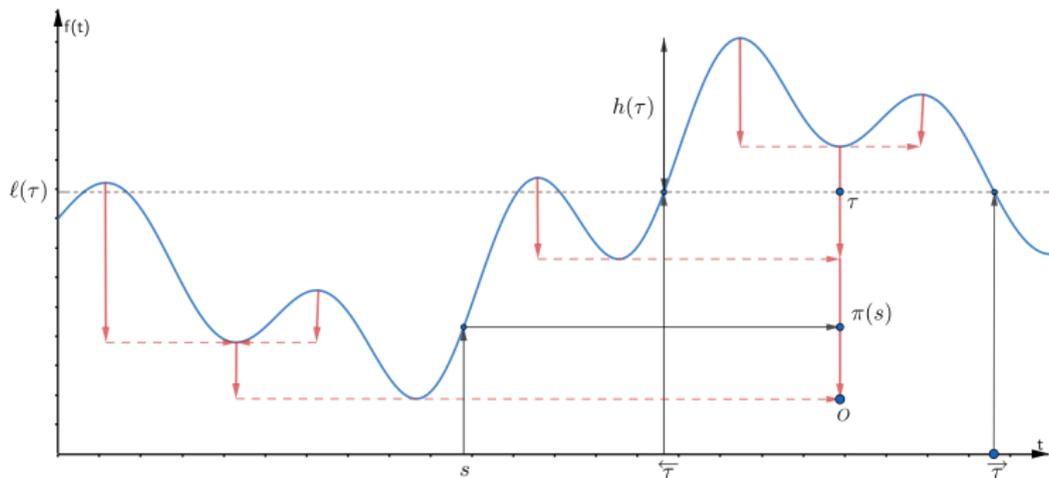


FIGURE – Arbre associé à une fonction sur  $[0, 1]$  et quantités associées.

## Notation

On notera  $\pi_f : X \rightarrow T_f$  la projection canonique, et  $h : T_f \rightarrow \mathbb{R}$  la fonction associant la distance entre le point  $\tau$  et la feuille la plus haute au-dessus de  $\tau$ .

## Définition

Avec la même fonction  $h$  sur l'arbre on peut définir l'arbre  $\varepsilon$ -simplifié

$$T_f^\varepsilon := \{\tau \in T_f \mid h(\tau) \geq \varepsilon\} \quad (1.3)$$

## Remarque

Par compacité et locale connexité de  $T_f$ ,  $T_f^\varepsilon$  est toujours un arbre fini.

## Définition

On définit  $N^\varepsilon$  comme étant le nombre de feuilles de l'arbre  $T_f^\varepsilon$ .

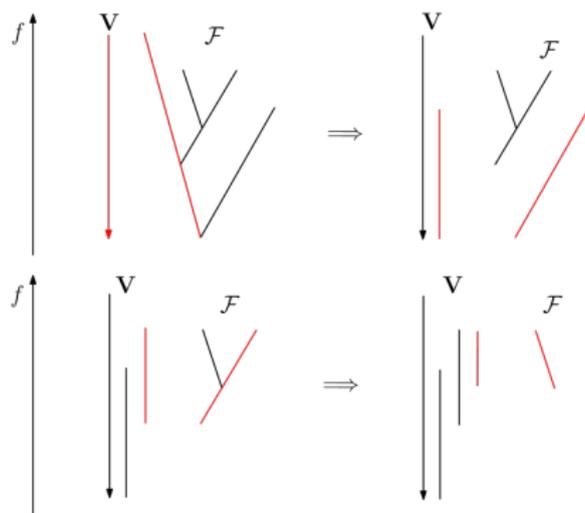


FIGURE – Construction du code-barres à partir de l'arbre.

## Remarque

*Si on fixe  $r$ , le code-barres quantifie le nombre de composantes connexes de  $\{f \geq r\}$ .  
En revanche, l'arbre contient plus d'informations que la donnée seule du code-barres.*

## Définition

Nous noterons  $N^\varepsilon$  le nombre de barres de  $\mathcal{B}(f)$  ayant une longueur  $\geq \varepsilon$ .

## Remarque

$N^\varepsilon$  correspond au nombre de variations de taille au moins  $\varepsilon$  de  $f$ .

## Définition

On définit la norme  $\ell_p$  d'une fonction  $f$  par

$$\ell_p(f) := \left[ \sum_{b \in \mathcal{B}(f)} \ell(b)^p \right]^{1/p}. \quad (2.4)$$

où  $\mathcal{B}(f)$  est le code-barres de  $f$ .

## Remarque

*Il est commode de faire le saut conceptuel de considérer la fonction  $\ell_p^p(f)$  pour  $p \in \mathbb{C}$ .*

## Proposition (P.)

*Les fonctions  $N^\varepsilon$  et  $\ell_p^p$  sont duales au sens où*

$$\ell_p^p(f) = p \int_0^\infty \varepsilon^{p-1} N^\varepsilon d\varepsilon \quad (2.5)$$

$$N^\varepsilon = \frac{1}{2\pi i} \int_{\gamma-i\infty}^{\gamma+i\infty} \varepsilon^{-p} \ell_p^p(f) \frac{dp}{p} \quad (2.6)$$

## Démonstration.

On exprime

$$\ell(b)^p = p \int_0^{\ell(b)} \varepsilon^{p-1} d\varepsilon = p \int_0^\infty 1_{[0, \ell(b)]}(\varepsilon) \varepsilon^{p-1} d\varepsilon, \quad (2.7)$$

en appliquant le théorème de Tonelli,

$$\ell_p^p(f) = p \int_0^\infty \underbrace{\left[ \sum_{b \in \mathcal{B}(f)} 1_{[0, \ell(b)]}(\varepsilon) \right]}_{=: N^\varepsilon} \varepsilon^{p-1} d\varepsilon. \quad (2.8)$$

Par le théorème d'inversion de Mellin, on a la deuxième équation. □

## Théorème (Picard, §3[8] et[6])

Pour une fonction continue  $f : [0, 1] \rightarrow \mathbb{R}$ ,

$$\mathcal{V}(f) = \mathcal{L}(f) = \limsup_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} \vee 1 \quad (2.9)$$

où  $a \vee b := \max\{a, b\}$  et

$$\mathcal{V}(f) := \inf\{p \mid \|f\|_{p\text{-var}} < \infty\} \quad \text{et} \quad \mathcal{L}(f) := \inf\{p \mid \ell_p < \infty\}. \quad (2.10)$$

## Remarque

*A priori*  $\ell_p^p$  est holomorphe sur le demi-plan  $\operatorname{Re}(p) > \mathcal{L}(f)$ .

Maintenant, considérons  $f$  un processus stochastique (p.s. continu) sur  $X$ , i.e.

## Définition

$f : \Omega \times X \rightarrow \mathbb{R}$  avec  $(\Omega, \mathcal{F}, \mathbb{P})$  un espace probabilisé.

## Définition

La fonction  $\zeta$  de  $f$  est définie par

$$\zeta_f(p) := \mathbb{E}[\ell_p^p(f)] \quad (3.11)$$

## Remarque

Pour des processus "canoniques",  $\zeta_f$  admet en fait un prolongement méromorphe à tout le plan complexe.

## Remarque

Dorénavant  $X = [0, t]$  et on autorisera  $f$  p.s. càdlàg.

## Fait (Flajolet [4])

*On a la correspondance suivante*

*Extensions méro. de  $\zeta$  (à gauche, à droite)*



*Développements (hyper)asymptotiques de  $\mathbb{E}[N^\varepsilon]$  (lorsque  $\varepsilon \rightarrow 0, \varepsilon \rightarrow \infty$ ).*

## Remarque

*Le fait ci-haut s'applique de manière générale aux fonctions et leur transformées de Mellin.*

## Théorème (P.)

La fonction  $\zeta$  du mouvement brownien  $B$  sur l'intervalle  $[0, t]$  admet un prolongement méromorphe partout sur  $\mathbb{C}$ . De plus, elle est exactement égale à

$$\zeta_B(p) = \frac{4(2^p - 3)}{\sqrt{\pi}} \left(\frac{t}{2}\right)^{\frac{p}{2}} \Gamma\left(\frac{p+1}{2}\right) \zeta(p-1) \quad (3.12)$$

pour tout  $p$  et elle admet un unique pôle simple en  $p = 2$  de résidu  $[B]_t = t$ .

## Théorème (P.)

En fait, la fonction  $\zeta$  de toute semimartingale  $S$  admet un pôle simple à  $p = 2$  de résidu  $[S]_t$ .

## Proposition (P.)

Dans le cas du mouvement brownien,  $\mathbb{E}[N^\varepsilon]$  admet les deux représentations suivantes qui convergent bien pour grand et petit  $\varepsilon$  respectivement

$$\mathbb{E}[N^\varepsilon] = \frac{t}{2\varepsilon^2} + \frac{2}{3} + o(\varepsilon^n) \text{ lorsque } \varepsilon \rightarrow 0. \quad (3.13)$$

## Démonstration.

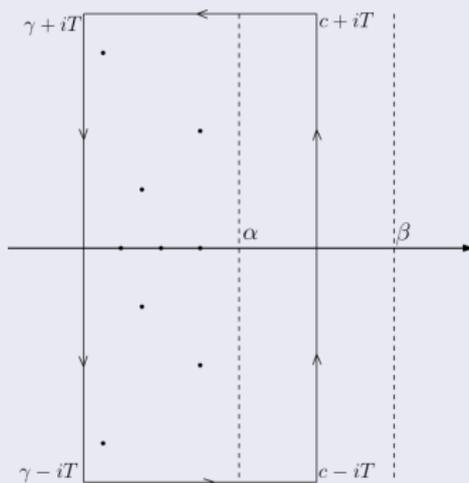
C'est la correspondance fondamentale. On applique l'inversion de Mellin à  $\zeta_B(p)/p$ .

$$\mathbb{E}[N^\varepsilon] = \frac{1}{2\pi i} \int_{\gamma-i\infty}^{\gamma+i\infty} \varepsilon^{-p} \zeta_B(p) \frac{dp}{p} \quad (3.14)$$

On a deux pôles, un à  $p = 2$  et un autre à  $p = 0$ . □

Suite.

On intègre le long du contour suivant

Par le théorème des résidus, on a le résultat. □

En fait on a mieux,

## Proposition (P.)

Dans le cas du mouvement brownien,  $\mathbb{E}[N^\varepsilon]$  admet les deux représentations suivantes qui convergent bien pour grand et petit  $\varepsilon$  respectivement

$$\mathbb{E}[N^\varepsilon] = 4 \sum_{k \geq 1} (2k - 1) \operatorname{erfc} \left( \frac{(2k - 1)\varepsilon}{\sqrt{2t}} \right) - k \operatorname{erfc} \left( \frac{2k\varepsilon}{\sqrt{2t}} \right) \quad (3.15)$$

$$= \frac{t}{2\varepsilon^2} + \frac{2}{3} + 2 \sum_{k \geq 1} (2(-1)^k - 1) \frac{e^{-\frac{\pi^2 k^2 t}{2\varepsilon^2}}}{\varepsilon^2} \left[ 1 + \frac{\varepsilon^2}{\pi^2 k^2 t} \right]. \quad (3.16)$$

## Remarque

On peut le démontrer en prenant la transformée de Mellin inverse explicitement. On obtient la deuxième expression comme image inverse par la transformée de Mellin de  $\zeta_B$  en utilisant l'équation fonctionnelle de la fonction  $\zeta$  de Riemann.

## Définition

Un processus stochastique  $f$  sur  $\mathbb{R}_+$  est dit de Lévy si

- 1 p.s.  $f(0) = 0$  ;
- 2 Les accroissements de  $f$  sont indépendants, i.e. pour toute partition finie  $(t_i)_i$  de  $\mathbb{R}_+$ ,  $(f(t_i) - f(t_{i-1}))_i$  sont mutuellement indépendants ;
- 3 Les accroissements de  $f$  sont stationnaires, i.e. pour tout  $s < t \in \mathbb{R}_+$ ,  $f(t - s) = f(t) - f(s)$  en distribution ;
- 4  $f$  est continu en probabilité.

## Définition

Un processus de Lévy est dit  $\alpha$  stable s'il existe  $\alpha$  tel que pour tout  $t$  et tout  $c$ ,  $f(c^\alpha t) = cf(t)$  en distribution.

## Remarque

En particulier, le mouvement brownien est un processus de Lévy 2-stable.

## Théorème (P.)

*Si  $f$  est un processus de Lévy  $\alpha$ -stable, il existe une variable aléatoire  $U$  avec des moments finis telle que la fonction  $\zeta$  du processus admet un prolongement méromorphe partout sur  $\mathbb{C}$  avec un seul pôle simple en  $p = \alpha$  de résidu  $\mathbb{E}[U] \alpha t$ .*

## Théorème (P.)

*Pour un processus de Lévy  $\alpha$ -stable sur  $[0, t]$ , on a le développement asymptotique suivant*

$$\mathbb{E}[N^\varepsilon] = \frac{t}{\mathbb{E}[U] \varepsilon^\alpha} + \frac{\mathbb{E}[U^2]}{2\mathbb{E}[U]^2} + o(\varepsilon^{\alpha n}) \quad \text{lorsque } \varepsilon \rightarrow 0, \quad (3.17)$$

*pour tout  $n \in \mathbb{N}$ . De plus, p.s.*

$$N^\varepsilon \sim \frac{t}{\mathbb{E}[U] \varepsilon^\alpha} \quad \text{lorsque } \varepsilon \rightarrow 0. \quad (3.18)$$

## Théorème (P., Théorème 3.1 [7])

Soit  $\delta_N := \|X - X_N\|_{L^\infty} \rightarrow 0$  p.s. quand  $N \rightarrow \infty$ , alors pour tout  $\varepsilon \geq 2\delta_N$

$$N_X^{\varepsilon+\delta_N} \leq N_{X_N}^\varepsilon \leq N_X^{\varepsilon-\delta_N} \quad (3.19)$$

De plus, si  $\mathbb{E}[N_X^\varepsilon]$  est continue en  $\varepsilon$ , alors

$$N_{X_N}^\varepsilon \xrightarrow[N \rightarrow \infty]{L^1} N_X^\varepsilon \quad \text{et} \quad N_{X_N}^\varepsilon \xrightarrow[N \rightarrow \infty]{\mathbb{P}} N_X^\varepsilon,$$

à  $N$  fixé, ceci se traduit quantitativement par

$$\mathbb{E}\left[|N_X^\varepsilon - N_{X_N}^\varepsilon|\right] \leq 2\omega_\varepsilon(\delta_N) \quad \text{et} \quad \mathbb{P}\left(|N_X^\varepsilon - N_{X_N}^\varepsilon| \geq k\right) \leq \frac{2\omega_\varepsilon(\delta_N)}{k} \quad (3.20)$$

où  $\omega_\varepsilon$  est le module de continuité de  $\mathbb{E}[N_X^\varepsilon]$  sur l'intervalle  $[\varepsilon - \delta_N, \varepsilon + \delta_N]$ . Enfin, on a les inégalités suivantes

$$N_{X_N}^{\delta_N} \geq N_X^{2\delta_N} \quad \text{and} \quad N_X^{\delta_N} \geq N_{X_N}^{2\delta_N}.$$

Soit  $f$  un processus  $\alpha$ -stable.

- ① Échantillonner  $M$  chemins du processus stochastique  $f$  (par exemple à intervalles réguliers de taille  $\frac{1}{N}$  pour un certain  $N$ );
- ② Calculer le code-barres des chemins échantillonnés.
- ③ Pour *une certaine plage de valeurs suffisamment petites*  $\varepsilon$ , et pour une un  $c > 1$ , calculer

$$\hat{\alpha}_M := \log_c \left[ \frac{\overline{N}_t^{\varepsilon/c} - \overline{N}_t^{2\varepsilon/c}}{\overline{N}_t^\varepsilon - \overline{N}_t^{2\varepsilon}} \right]. \quad (3.21)$$

Ici, la notion de *une certaine plage de valeurs suffisamment petites*  $\varepsilon$  et la constante  $c$  dépendent toutes deux de  $N$ , avec la condition limite que lorsque  $N \rightarrow \infty$ , la borne inférieure de la plage de  $\varepsilon$  valides s'approche de 0.

En ignorant les contributions superpolynomiales l'argument à l'intérieur du log de l'estimateur est approx.

$${}_c \hat{\alpha}_M \approx \frac{\frac{t}{\mathbb{H}[U](\varepsilon/c)^\alpha} - \frac{t}{\mathbb{H}[U](2\varepsilon/c)^\alpha}}{\frac{t}{\mathbb{H}[U]\varepsilon^\alpha} - \frac{t}{\mathbb{H}[U](2\varepsilon)^\alpha}} \approx c^\alpha. \quad (3.22)$$

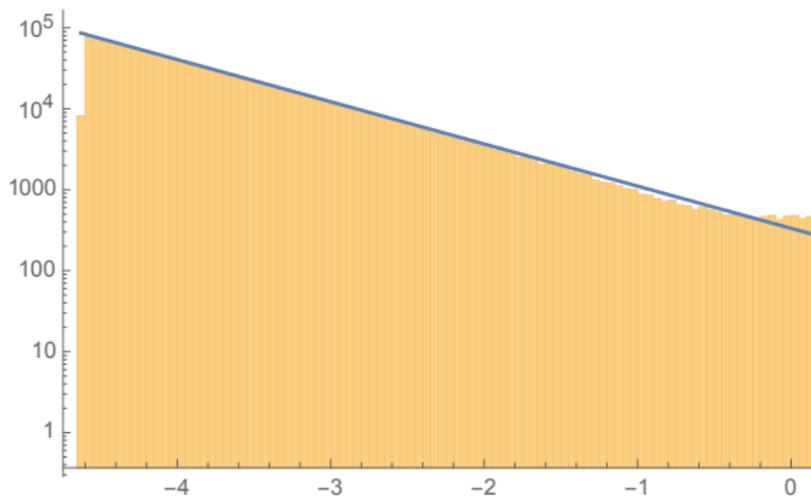


FIGURE – Test statistique pour un processus 1.2-stable. La pente de la droite est bien 1.2.

- Les prolongements analytiques de  $\zeta_f$  ne sont pas triviaux. Il serait intéressant d'étudier les séries de Dirichlet associées (déjà quelques résultats dans cette direction...);
- $\ell_p$  n'est pas stable par perturbations  $L^\infty$ , mais sa variable duale  $N^\varepsilon$  l'est et semble être robuste statistiquement au vu de nos résultats pour les processus  $\alpha$ -stables;
- Ceci donne une explication qualitative sur le fait que  $\ell_p$  soit mieux adapté dans le cadre de l'apprentissage statistique;
- Il serait intéressant d'étudier les fonctions  $\zeta$  de fonctions sur des espaces plus généraux de plus grande dimension.

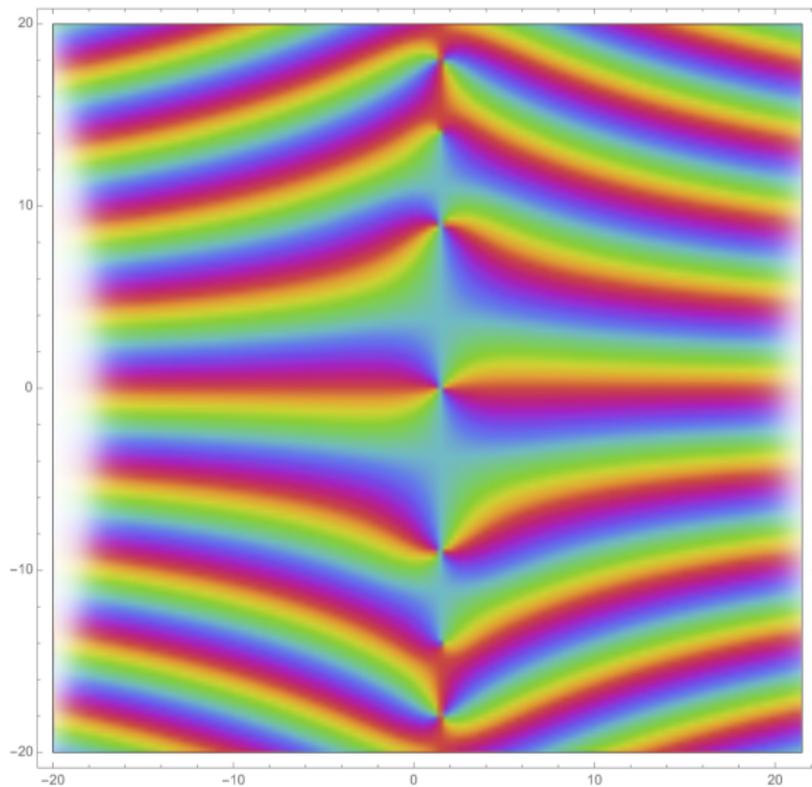


FIGURE – La fonction  $\eta_B$ .

## Notation

Dans la suite, soit  $X$  une variété et  $f : X \rightarrow \mathbb{R}$  une fonction continue.

## Remarque

La fonction  $f$  induit une filtration de l'espace  $X$  par

$$X_r := \{f \geq r\} \tag{4.23}$$

## Définition

Fixons un corps  $k$ . L'homologie persistante du couple  $(X, f)$  est un foncteur de la petite catégorie (partiellement ordonnée)  $\mathbb{R}$  vers  $\text{Vect}_k$

$$H_*(X, f) : r \mapsto H_*(X_r, k) \tag{4.24}$$

$$[s \rightarrow r] \mapsto [H_*(X_s, k) \rightarrow H_*(X_r, k)] \tag{4.25}$$

## Théorème (Décomposition, [2, 5])

Pour les fonctions continues  $f : X \rightarrow \mathbb{R}$ , l'homologie persistante admet toujours une décomposition en des "modules d'intervalle"

$$H_*(X, f) := \bigoplus_{(a,b)} k[a, b[ \quad (4.26)$$

où les modules d'intervalle sont des foncteurs  $\mathbb{R} \rightarrow \text{Vect}_k$  définis par :

$$k[a, b[ : r \mapsto \begin{cases} k & \text{si } r \in [a, b[ \\ 0 & \text{sinon} \end{cases} \quad (4.27)$$

$$[s \rightarrow r] \mapsto \begin{cases} \text{id} & \text{si } s, r \in [a, b[ \\ 0 & \text{sinon} \end{cases} \quad (4.28)$$

## Remarque

Toute l'information de  $H_*(X, f)$  est contenue dans les couples  $(a, b) \in \mathbb{R}^2$ .

## Remarque

*Dans la suite, nous ne considérerons l'homologie persistante qu'en degré 0. Dans ce cas, on s'intéresse juste aux composantes connexes.*

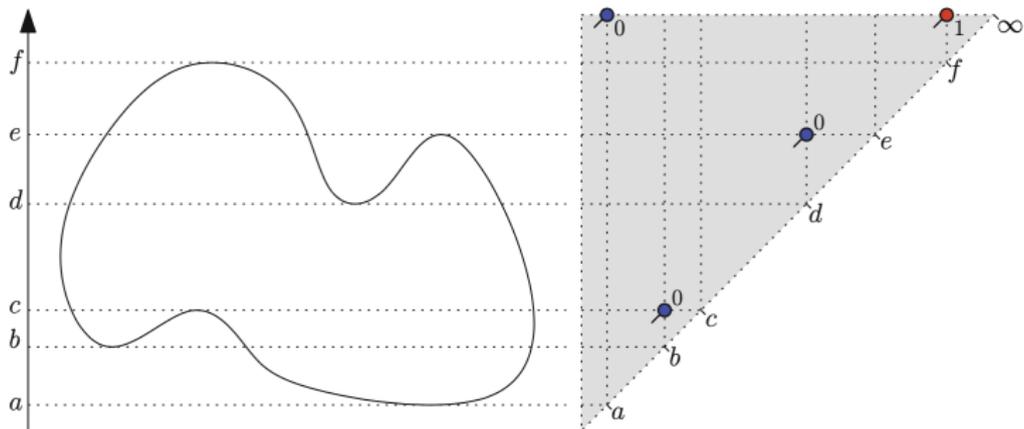


FIGURE – Diagramme de persistance associé à une filtration du cercle par ensembles de sous-niveaux [2].

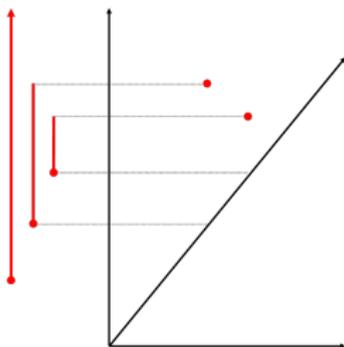


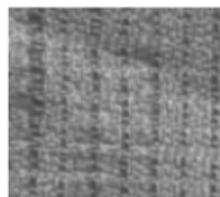
FIGURE – Diagrammes de persistance et code-barres.

## Remarque

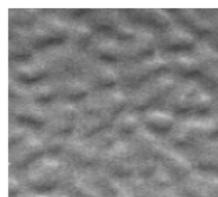
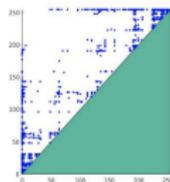
*Il existe une distance naturelle sur l'espace des diagrammes appelée "distance bottleneck", stable par perturbations  $C^0$  de  $f$ . Cette distance "oublie" les petites barres du code-barres.*

## Remarque

*En pratique, les petites barres identifient les “phénomènes haute fréquence” comme la texture. Pour certaines applications, elles sont donc fondamentales.*



Label = Canvas



Label = Carpet

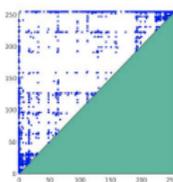


FIGURE – Diagrammes associés à différentes textures [1].

## Définition

$$d_p(a, b) = \left[ \inf_{\pi \in \Gamma(a, b)} \sum_{x \in a \cup \partial\Omega} d(x, \pi(x))^p \right]^{\frac{1}{p}} \quad (4.29)$$

où  $d$  est la norme  $q$  sur  $\mathbb{R}^2$ ,  $\Gamma(a, b)$  est l'ensemble des bijections entre  $a \cup \partial\Omega$  et  $b \cup \partial\Omega$ .

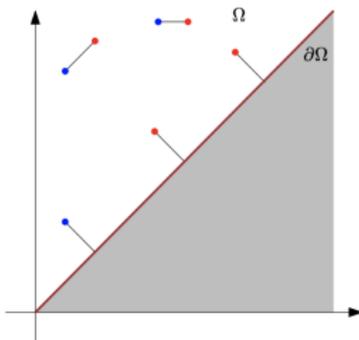


FIGURE – Matching partiel entre deux diagrammes  $a$  et  $b$  [3].

## Remarque

$\ell_p$  est la distance entre le diagramme et la diagonale.

## Théorème

Si  $f = M + A$  est une semimartingale continue tel que pour tout  $s \geq 1$

$$\mathbb{E} \left[ [M]_t^s + \left( \int_0^t |dA|_s \right)^s \right] < \infty. \quad (4.30)$$

Alors dans  $L^s$  on a

$$N^\varepsilon \sim \frac{[f]_t}{2\varepsilon^2} \quad \text{lorsque } \varepsilon \rightarrow 0. \quad (4.31)$$

En particulier  $\zeta_f$  admet un pôle simple à  $p = 2$  de résidu  $[f]_t$ .

## Remarque

Si  $f = B$ , l'énoncé du théorème est vrai p.s.

 M. Carriere, S. Oudot, and M. Ovsjanikov.  
Sliced Wasserstein Kernel for Persistence Diagrams.  
*In ICML 2017 - Thirty-fourth International Conference on Machine Learning*,  
pages 1–10, Sydney, Australia, Aug. 2017.

 F. Chazal, V. de Silva, M. Glisse, and S. Oudot.  
*The Structure and Stability of Persistence Modules*.  
Springer International Publishing, 2016.

 V. Divol and T. Lacombe.  
Understanding the topology and the geometry of the persistence diagram space  
via optimal partial transport.  
*ArXiv*, abs/1901.03048, 2019.

 P. Flajolet, X. Gourdon, and P. Dumas.  
Mellin transforms and asymptotics : Harmonic sums.  
*Theoretical Computer Science*, 144(1) :3–58, 1995.

 S. Y. Oudot.  
*Persistence Theory - From Quiver Representations to Data Analysis*, volume 209  
of *Mathematical surveys and monographs*.  
American Mathematical Society, 2015.

 D. Perez.  
On  $C^0$ -persistent homology and trees.  
<https://arxiv.org/abs/2012.02634>, Dec. 2020.

 D. Perez.

On the persistent homology of almost surely  $C^0$  stochastic processes.  
<https://arxiv.org/abs/2012.09459>, Dec. 2020.



J. Picard.

A tree approach to  $p$ -variation and to integration.

*The Annals of Probability*, 36(6) :2235–2279, Nov 2008.