

# Topological Data Analysis

D'après Julien Tierny et Frédéric Chazal

4 novembre 2025



## Table des matières

<b>1 Homologie Simpliciale</b>	<b>2</b>
1.1 Complexe Simpliciaux . . . . .	2
1.2 Homologie simpliciale . . . . .	4
1.2.1 Rappels d'homotopie . . . . .	4
1.2.2 Groupes d'homologie . . . . .	4
<b>2 Homologie Persistente</b>	<b>5</b>
2.1 Invariance Topologique et Filtrations . . . . .	5
2.1.1 Homologie Singulière . . . . .	5
2.1.2 Filtrations . . . . .	6
2.2 Homologie Persistente . . . . .	7
2.2.1 Sur les fonctions . . . . .	7
2.2.2 Sur les filtrations . . . . .	7
2.3 Calcul de filtrations . . . . .	9
2.3.1 Complexes de Čech, de Vietoris-Rips, et autres . . . . .	9
2.3.2 Stabilité . . . . .	10
2.3.3 Théorèmes de Stabilité . . . . .	11
2.4 Calculabilité et bruit . . . . .	13
2.4.1 Calcul statistique . . . . .	13
2.4.2 Paysages de persistence et rééchantillonnage . . . . .	14
2.4.3 Bruit et méthode d'échantillonnage . . . . .	15
2.5 Persistence et apprentissage automatique . . . . .	16
2.5.1 Représentation de persistence . . . . .	16
2.5.2 Différentiabilité de la persistence . . . . .	17
<b>3 Fonctions de Morse</b>	<b>17</b>
<b>4 Inférence Topologique</b>	<b>17</b>
<b>5 Théorie de Morse Discrète</b>	<b>17</b>

### Résumé

<mailto:julien.tierny@sorbonne-universite.fr> <mailto:frederic.chazal@inria.fr> <https://julien-tierny.github.io/topologicalDataAnalysisClass.html> <https://geometrica.saclay.inria.fr/team/Fred.Chazal/MVA2025>

## Introduction

Méthodes algorithmiques d'analyse topologique de données, particulièrement en science et en ingénierie.

Le but est de partir de données, sous forme de maillages et maillables, et de retrouver des structures au sein de jeux de données. Partant d'une carte (considérée comme jeu de données brutes), avec des features intéressantes, pour pouvoir raisonner sur l'espace, on passe à une représentation abstraite, par exemple comme un graphe, et c'est sur cette structure de données sous-jacente qu'on va raisonner. Ici, on peut ajouter des filtres pour redéfinir le maillage et donc redéfinir le résultat du raisonnement. Plus généralement, on veut construire une carte à partir d'un jeu de données. En astrophysique, par exemple, on modélise la croissance de l'univers à une grande échelle, on la simule par une grille de voxel, on estime la densité de matière noire sur chaque voxel, et on découvre une sorte de géométrie ressemblant à des neurones lorsqu'on trouve aussi des groupes de galaxies, formant une "toile cosmique". On peut calculer les connexions avec des complexes simpliciaux dits de Morse-Smale, dont on peut extraire une structure de graphes.

Ainsi, on extrait de la structure d'un ensemble de données, de manière robuste et indépendante de l'échelle, par comparaison et extraction de propriétés. Sous le capot, on fait :

- Homologie Simpliciale
- Théorie de Morse
- Homologie Persistente

Pour des données numériques, étant données un échantillon de points dans un espace euclidien, par exemple, on peut les représenter et objectiver des représentations géométriques apparaissant. On a des manières de mailler l'ensemble (triangulation de Delaunay, par exemple) qui amènent à des indicateurs qui nous expliquent où sont répartis les données, par exemple avec des noyaux pour estimer la densité. Avec une fonction scalaire sur un maillage, on définit une filtration, et on regarde les propriétés de la fonction, comme les optima locaux et on en extrait une structure algébrique (complexe de Morse-Smale) qui nous donne une structure algébrique. On obtient des générateurs, et des composantes "connexes".

On a ce genre de densité de pixels, par exemple la hauteur de surface de la mer qui permet de remarquer les vortexs, en chimie quantique ou des spectrogrammes d'enregistrement vocaux. On part d'un domaine géométrique et d'un signal sur ce domaine, signal qui exhibe des patterns géométriques qu'on souhaite quantifier. Ceci permet l'extraction de propriétés, la segmentation, la réduction de dimension et autres. Dans le cas de points en grande dimension, on a une unique théorie qui s'applique très généralement.

En terme de logiciels, on a le TTK (ParaView >= 5.10) et Gudhi (bibliothèque python).

## 1 Homologie Simpliciale

Les données reçues, parfois, vont contenir explicitement la géométrie avec une construction combinatoire. On supposera qu'on aura une donnée d'entrée linéaire par morceau sur un complexe simplicial.

### 1.1 Complexes Simpliciaux

**Définition 1.1** Un *d-simplexe* est l'enveloppe convexe  $\sigma$  de  $d+1$  points affinement indépendants dans l'espace euclidien  $\mathbb{R}^n$  avec  $0 \leq d \leq n$ . On dit que  $d$  est la *dimension* du simplexe.

Une ligne est un 1-simplexe, un triangle un 2-simplexe et un tétraèdre un 3-simplexe.

**Définition 1.2** Une *face*  $\tau$  d'un simplexe  $\sigma$  est un simplexe construit par un ensemble non vide des  $d+1$  points définissant  $\sigma$ . On note  $\tau \leq \sigma$  et  $\tau_i$  une face de dimension  $i$ . On dit aussi que  $\sigma$  est une *coface* de  $\tau$ .

Selon la définition, on a  $\sigma \leq \sigma$ .

**Définition 1.3** Un *complexe simplicial*  $\mathcal{K}$  est une collection finie non-vide de simplexes  $\{\sigma_i\}$ , telle que :

1.  $\tau \leq \sigma \Rightarrow \tau \in \mathcal{K}$  ;
2.  $\sigma_i \cap \sigma_j$  est soit une face, soit vide.

**Définition 1.4** *L'étoile* d'un simplexe  $\sigma \in \mathcal{K}$  est l'ensemble des simplexes de  $\mathcal{K}$  qui contiennent  $\sigma$  :

$$\text{St}(\sigma) = \{\tau \in \mathcal{K} | \sigma \leq \tau\}$$

On note  $\text{St}_d(\sigma)$  les  $d$ -simplexes de  $\text{St}(\sigma)$ .

C'est l'ensemble des cofaces de  $\sigma$  dans  $\mathcal{K}$ . C'est le plus petit voisinage combinatoire autour d'un simplexe.

**Définition 1.5** Le *lien* d'un simplexe  $\sigma$  est l'ensemble des faces de  $\text{St}(\sigma)$  disjointes de  $\sigma$  :

$$\text{Lk}(\sigma) = \{\tau \leq \sigma' | \sigma' \in \text{St}(\sigma) \wedge \tau \cap \sigma = \emptyset\}$$

On définit de même le *d-lien*  $\text{Lk}_d(\sigma)$  en remplaçant  $\text{St}$  par  $\text{St}_d$  dans la définition

C'est en quelque sorte la bordure du voisinage combinatoire de lui-même.

En réalité on va considérer que les sommets (ou 0-simplexes) sont des points, et que les  $d$ -simplexes sont des ensembles de points. Ceci définit une notion de complexe simplicial abstrait, utile lorsqu'on n'a pas d'immersion dans un espace euclidien, ou alors dans un espace euclidien en trop grande dimension. On relâche ici la condition d'intersection, puisqu'on n'a plus de structure géométrique de l'espace.

Un exemple de complexe simplicial abstrait est le complexe de Rips ou complexe de Vietori-Rips. Étant donné un nuage de points avec une métrique :

- Le diamètre d'un ensemble  $P$  est défini par :  $\emptyset(P) = \sup\{d(x, y) | x, y \in P\}$
- On construit un complexe simplicial  $p \leq p_{max}$  de sorte que tous  $p + 1$  points dont le diamètre est plus petit qu'une valeur seuil  $d_{max}$ .

Le complexe de Rips est une généralisation de la notion de graphe de voisinage.

**Définition 1.6** *L'espace sous-jacent* à un complexe simplicial est l'union des simplexes du complexe.

**Définition 1.7** La *triangulation*  $\mathcal{T}$  d'un espace topologique  $X$  est un complexe simplicial  $\mathcal{K}$  dont l'espace sous-jacent  $|\mathcal{K}|$  est homéomorphe à  $X$ .

Une triangulation d'un espace est donc un complexe simplicial abstrait.

**Définition 1.8** Une *d-variété*  $M$  est un espace topologique dans lequel tout élément  $m$  a un voisinage ouvert homéomorphe à une  $d$ -boule euclidienne.

**Définition 1.9** La triangulation d'une  $d$ -variété est appelée *d-variété linéaire par morceaux*.

Représenter en mémoire un complexe simplicial est très couteux : il faut, pour chaque dimension, une liste des hyperarêtes (ou  $d$ -simplexes) en les représentant par un indice de sommet.

## 1.2 Homologie simpliciale

### 1.2.1 Rappels d'homotopie

**Définition 1.10** Un *chemin*  $p$  dans  $C$  est un homéomorphisme d'un intervalle réel vers l'objet  $C$ . On dit que  $C$  est *connexe (par arcs)* si pour tous deux points il existe un chemin dans  $C$  les reliant.

**Définition 1.11** Une *composante connexe* d'un objet est un sous-ensemble connexe (par arcs) maximal de l'objet.

**Définition 1.12** Une *homotopie* entre deux fonctions continues  $f$  et  $g$  de  $X$  vers  $Y$  est une fonction continue  $H : X \rightarrow [0, 1] \rightarrow Y$  du produit d'un espace topologique  $X$  par l'intervalle unité vers un espace topologique  $Y$  de sorte que  $H(x, 0) = f(x)$  et  $H(x, 1) = g(x)$  pour tout  $x \in X$ . S'il existe une homotopie entre  $f$  et  $g$  on dit que  $f$  et  $g$  sont *homotopes*.

**Définition 1.13** Si dans un espace  $X$ , tous les chemins entre tous deux points sont homotopes, on dit que  $X$  est *simplement connexe*.

Le disque est simplement connexe, mais pas le disque privé de 0.

**Définition 1.14** La *caractéristique d'Euler* d'une triangulation  $T$  d'un espace topologique est la somme alternée des nombres des  $i$ -simplexes :

$$\chi(T) = \sum_{i=0}^d (-1)^i |\sigma_i|$$

**Proposition 1.1** La caractéristique d'Euler est invariante par homéomorphisme.

### 1.2.2 Groupes d'homologie

**Définition 1.15** Une  *$p$ -chaîne* est une somme (formelle) de  $p$ -simplexes. On suppose que l'opérateur somme est défini modulo 2.

Ici, la somme est réellement la différence symétrique (ou la disjonction exclusive) sur  $\mathbb{F}_2^{|\sigma_p|}$ , et définit le groupe  $C_p$  des  $p$ -chaînes. Le rang de  $C_p$  est  $|\sigma_p|$  et son ordre est  $2^{|\sigma_p|}$ .

On peut généraliser à des coefficients plus généraux. Informatiquement, il faut voir la notion de chaîne comme un masque binaire sur l'ensemble des  $p$ -simplexes, et l'addition comme la disjonction exclusive.

**Définition 1.16** *L'opérateur de bordure*  $\partial$  d'un  $p$ -simplexe renvoie la  $(p - 1)$ -chaîne des  $(p - 1)$ -faces du simplexe. On l'étend aux  $p$ -chaînes comme un morphisme de  $C_p \rightarrow C_{p-1}$ .

Pour un triangle, c'est l'ensemble de ses arêtes. Pour une arête, c'est l'ensemble des extrémités.

**Proposition 1.2** L'opérateur de bordure définit une suite exacte appelée le complexe de chaîne associé au complexe  $K$  de dimension  $d$  :

$$\{0\} \rightarrow C_d(K) \xrightarrow{\partial} C_{d-1}(K) \xrightarrow{\partial} \cdots \xrightarrow{\partial} C_{k+1}(K) \xrightarrow{\partial} C_k(K) \xrightarrow{\partial} \cdots \xrightarrow{\partial} C_1(K) \xrightarrow{\partial} C_0(K) \xrightarrow{\partial} 0$$

**Définition 1.17** Un  *$p$ -cycle* est une  $p$ -chaîne dont la bordure est vide. On définit  $Z_p$  le groupe des  $p$ -cycles comme sous-groupe de  $C_p$ .

**Définition 1.18** Le groupe  $B_p$  des  $p$ -bordures est l'image de  $C_{(p+1)}$  par  $\partial$ .

**Lemme 1.1** Pour tout  $x \in C_p, p \geq 2, \partial\partial x = 0$ .

*Démonstration.* Il suffit de vérifier le résultat sur les  $p$ -simplexes et d'étendre par somme. Puisque pour tout  $(p-2)$ -faces  $\tau$  on a exactement 2  $(p-1)$ -co-faces de  $\tau$  dans un  $p$ -simplexe  $\sigma$ , on a le résultat. ■

On obtient directement :

**Proposition 1.3**  $B_p$  est un sous-groupe de  $Z_p$ .

Si on a trois 1-simplexes  $e_1, e_2, e_3$  mais pas leur coface commune  $\{e_1, e_2, e_3\}$ , on a un exemple d'inclusion stricte. Ceci nous amène à définir la notion de trou, en isolant les cycles :

**Définition 1.19** Le  $p$ -ème groupe d'homologie  $H_p$  est le quotient de  $Z_p$  par  $B_p$ .

*Démonstration.*  $B_p$  étant un sous-groupe de  $Z_p$ ,  $H_p$  est bien défini. ■

Géométriquement, on peut étendre un  $p$ -cycle à un autre  $p$ -cycle lorsqu'ils encapsulent le même "trou", c'est-à-dire lorsque qu'on peut "étendre" le premier cycle en encapsulant un  $(p+1)$ -simplexe. Une classe d'homologie est un élément de  $H_p$ , ou plutôt sa classe d'équivalence dans  $Z_p$ .

Pour calculer  $|H_p|$ , on énumère  $C_p$ , on élimine les chaînes de bordure non-vide pour calculer  $Z_p$ , et on peut ensuite énumérer les classes d'homologie.

**Définition 1.20** On définit le  $p$ -ème nombre de Betti  $\beta_p$  comme le rang du groupe  $H_p$ . Ici, c'est  $\log_2 |H_p|$ .

La formule logarithmique pour  $\beta_p$  vient du calcul modulo 2 dans notre opération de groupe.

**Proposition 1.4** La caractéristique d'Euler d'une triangulation  $T$  d'un espace topologique  $X$  de dimension  $d$  vérifie :

$$\chi(T) = \sum_{i=0}^d (-1)^i \beta_i(T)$$

On a des interprétations des nombres de Betti en faible dimension. Par exemple,  $\beta_0(K)$  est le nombre de composantes connexes de  $K$ .

## 2 Homologie Persistente

### 2.1 Invariance Topologique et Filtrations

#### 2.1.1 Homologie Singulière

On a le résultat important suivant, qui va nous permettre de simplifier la manière de définir l'homologie, notamment dans la représentation des simplexes :

**Théorème 2.1** Si  $K$  et  $K'$  sont des complexes simpliciaux de supports homéomorphes, leurs groupes d'homologie sont isomorphes et leurs nombres de Betti sont égaux.

*Démonstration.* « Débrouille-toi, normalien. » ■

On va donc définir l'homologie singulière, sur tout espace topologique, en considérant à homotopie près. On note  $\Delta_k$  le  $k$ -simplexe standard dans  $\mathbb{R}^k$ .

**Définition 2.1** Un *k-simplexe singulier* dans un espace topologique  $X$  est une fonction continue  $\sigma : \Delta_k \rightarrow X$ .

On reprends les constructions de l'homologie simpliciale pour les complexes singuliers : c'est la définition de l'homologie singulière.

**Proposition 2.1** L'homologie singulière est définie pour tout espace  $X$ . Si  $X$  est homotopiquement équivalent au support d'un complexe simpliciel, les homologies singulières et simpliciales coïncident.

**Proposition 2.2** Si  $f : X \rightarrow Y$  est continue, et  $\sigma : \Delta_k \rightarrow X$  est un simplexe dans  $X$ , alors  $f \circ \sigma$  est un simplexe dans  $Y$ , et  $f$  définit donc une application linéaire entre les groupes d'homologie :

$$f_{\sharp} : H_k(X) \rightarrow H_k(Y)$$

Si  $f$  est un homéomorphisme ou une équivalence d'homotopie, alors  $f_{\sharp}$  est un isomorphisme.

**R** En particulier, si  $X \subset Y$ , l'application d'inclusion induit une application linéaire d'homologie.

### 2.1.2 Filtrations

**Définition 2.2** Un *complexe simplicial filtré (ou une filtration)*  $\mathbb{K}$  construit sur un ensemble  $X$  est une famille  $\{K_a \mid a \in T\}$ , où  $T \subseteq \mathbb{R}$ , de sous-complexes d'un certain complexe simplicial fixé  $K$  avec ensemble de sommets  $X$  de sorte que  $K_a \subseteq K_b$  pour tout  $a \leq b$ .

L'homologie persistente d'un complexe simplicial filtré encode l'évolution de l'homologie des sous-complexes.

**Définition 2.3** Une *filtration d'un complexe simplicial fini*  $K$  est une séquence de sous-complexes telle que :

1.  $\emptyset = K^0 \subset K^1 \subset \cdots \subset K^m = K$
2.  $K^{i+1} = K^i \cup \sigma^{i+1}$  où  $\sigma^{i+1}$  est un simplexe de  $K$ .

La famille des ensembles de sous-niveau pour une fonction est un exemple de filtration. On verra plus bas comment définir de telles filtrations dans plus de cas. On a toutefois un algorithme pour calculer itérativement l'homologie simpliciale d'un complexe étant donné une filtration de celui-ci :

---

#### Algorithme 1 Calcul d'homologie (simpliciale)

---

**Input** Une filtration  $(K^i)_{i \leq m}$  d'un complexe simplicial  $K$  en dimension  $d$

```

 $\beta_i \leftarrow 0$ 
for  $i \in \{1, \dots, m\}$  do
     $k \leftarrow \dim \sigma^i - 1$ 
    if  $\sigma^i$  est dans un  $(k + 1)$ -cycle de  $K^i$  then
         $\beta_{k+1} \leftarrow \beta_{k+1} + 1$ 
    else
         $\beta_k \leftarrow \beta_k - 1$ 
return  $(\beta_0, \dots, \beta_d)$ 

```

---

**Définition 2.4** Un  $(k + 1)$ -simplexe  $\sigma^i$  est dit *positif* s'il est contenu dans un  $(k + 1)$ -cycle de  $K^i$ . Il est dit *négatif* sinon.

Un  $(k + 1)$ -simplexe positif crée un  $(k + 1)$ -cycle dans  $K^i$ . Un  $(k + 1)$ -simplexe négatif détruit un  $k$ -cycle dans  $K^i$ .

On veut donc vérifier :

$$\beta_k(K) = |k\text{-simplexes positifs}| - |(k - 1)\text{-simplexes négatifs}|$$

On a :

**Lemme 2.1** Si  $\sigma^i$  est un  $k$ -simplexe positif, il existe un unique  $k$ -cycle  $c_\sigma$  tel que :

1.  $c_\sigma$  n'est pas une bordure dans  $K^i$
2.  $c_\sigma$  contient  $\sigma^i$  mais pas d'autre  $k$ -simplexe positif.

*Démonstration.* Par induction sur l'ordre d'apparition des simplexes dans la filtration. ■

*Preuve de correction de l'algorithme 1.* Si  $\sigma^i$  est contenu dans un  $(k+1)$ -cycle  $c$  de  $K^i$ , alors ce cycle n'est pas une bordure dans  $K^i$ . De plus,  $c$  ne peut pas être homologue à un cycle dans  $K^{i-1}$ , et donc  $\beta_{k+1}(K^i) \geq \beta_{k+1}(K^{i-1}) + 1$ . Si  $\sigma^i$  n'est contenu dans aucun  $(k+1)$ -cycle  $c$  de  $K^i$ , alors  $\partial\sigma^i$  n'est pas une bordure dans  $K^{i-1}$  et donc  $\beta_k(K^i) \leq \beta_k(K^{i-1}) - 1$ . ■

Cela pose quelques questions, qui vont nous conduire à introduire l'homologie persistente.

## 2.2 Homologie Persistente

### 2.2.1 Sur les fonctions

Pour définir les diagrammes de persistence pour une fonction  $f : X \rightarrow \mathbb{R}$ , on étudie ses ensembles de sous-niveau. On représente, pour chaque dimension  $d$  d'homologie, la "durée de vie" d'une propriété topologique de dimension  $d$ . Ces propriétés topologiques sont, entre autres, observées par la variation du  $d$ -ème nombre de Betti. On représente alors sur un graphe 2-D, en abscisse, la valeur  $x$  pour laquelle une propriété apparaît, et en ordonnée la valeur  $x'$  pour laquelle la propriété disparaît. On a nécessairement  $x' > x$ . On notera  $D_{f,d}$  le diagramme défini ci-dessus, comme son ensemble de points du plan  $\mathbb{R}_+^2$ .

On définit alors une distance sur deux tels diagrammes :

**Définition 2.5** La *distance infinie de Wasserstein, ou distance du goulot* entre deux diagrammes  $D_1$  et  $D_2$  est définie par :

$$d_B(D_1, D_2) = \inf_{\gamma \in \Gamma} \sup_{p \in D_1} \|p - \gamma(p)\|_\infty$$

où  $\Gamma$  est l'ensemble des bijections entre  $D_1$  et  $D_2$ .

On note que, les normes étant équivalentes sur  $\mathbb{R}^2$ , la distance de Wasserstein 2 est, à un facteur près, la distance infinie.

**R** Pour pouvoir obtenir des bijections, on doit souvent *augmenter* les diagrammes, en projetant les points de l'un sur la diagonale de l'autre et réciproquement.

On aura alors un théorème important de stabilité :

**Théorème 2.2** Pour toutes fonctions *dociles*  $f, g : \mathbb{X} \rightarrow \mathbb{R}$ ,  $d_B(D_f, D_g) \leq \|f - g\|_\infty$ .

Cette définition, avec les mains, va être précisée et étendue plus bas.

### 2.2.2 Sur les filtrations

On va maintenant étendre la notion de diagrammes de persistences aux filtrations de complexes simpliciaux. On a une relation fondamentale :

**Proposition 2.3** Si  $t \leq t'$ ,  $f^{-1}([-\infty, t]) \subseteq f^{-1}([-\infty, t'])$ . Si  $f$  est définie sur les sommets d'un complexe simplicial  $K$ , et étendue de sorte que

$$f(\sigma = [v_0, \dots, v_k]) = \max f(v_i),$$

alors les ensembles de sous-niveau de  $f$  définissent une filtration du complexe simplicial  $K$ .

Il suffit maintenant d'adapter l'algorithme 1 ci-dessus pour maintenir une base d'homologie et les paires naissance-mort d'une propriété. On notera  $H_k^i = H_k(K^i)$ , et on va construire des bases par récurrence (les  $H_k$  sont des  $\mathbb{F}_2$ -espaces vectoriels).

La base de  $H_k^0$  est vide, puisque l'ensemble est vide. Si on a construit une base de  $H_k^{i-1}$ , on a deux cas :

1. Si  $\sigma^i$  est un  $k$ -simplexe positif, alors on ajoute la classe d'homologie du cycle  $c^i$  associé à  $\sigma^i$  par le Lemme 2.1 à la base de  $H_k^{i-1}$  pour obtenir une base de  $H_k^i$ .
2. Si  $\sigma^i$  est un  $(k+1)$ -simplexe négatif :
  - On dénote  $c^{j_1}, \dots, c^{j_p}$  les cycles associés aux simplexes positifs  $\sigma^{j_1}, \dots, \sigma^{j_p}$  de la base de  $H_k^{i-1}$ .
  - On pose  $d = \partial\sigma^j = \sum_{k=1}^p \varepsilon_k c^{j_k} + b$
  - On pose  $l(i) = \max \{j_k \mid \varepsilon_k = 1\}$
  - On enlève la classe d'homologie de  $c^{l(i)}$  pour obtenir une base de  $H_k^i$ .

Ceci explique comment modifier l'algorithme 1 pour calculer les diagrammes de persistence. Cependant, avant de réécrire l'algorithme, on va s'intéresser à un test algorithmique pour vérifier que  $\sigma^j$  est positif ou négatif. Pour ce faire, on introduit la matrice de l'opérateur de bordure. On rappelle qu'on se donne une filtration : d'un complexe simplicial fini  $d$ -dimensionnel  $\emptyset = K^0 \subset K^1 \subset \dots \subset K^m = K$  telle que  $K^{i+1} = K^i \cup \sigma^{i+1}$  où  $\sigma^{i+1}$  est un simplexe de  $K$ .

**Définition 2.6** On pose  $M = (m_{i,j})_{1 \leq i,j \leq m}$  telle que  $m_{i,j} = 1$  si, et seulement si  $\sigma^i$  est une face de  $\sigma^j$  et vaut 0 sinon. C'est la *matrice de l'opérateur d'inclusion*.

Pour toute colonne  $C_j$ , on définit donc  $l(j)$  par :

$$(i = l(j)) \Leftrightarrow (m_{i,j} = 1 \wedge m_{i',j} = 0, \forall i' > i)$$

On obtient une version matricielle de l'algorithme de persistence :

---

### Algorithme 2 Algorithme de Persistence, version Matricielle

---

**Input** Une filtration  $\emptyset = K^0 \subseteq \dots \subseteq K^m = K$  d'un complexe simplicial  $d$ -dimensionnel de sorte que  $K^{i+1} = K^i \cup \sigma^{i+1}$  où  $\sigma^{i+1}$  est un simplexe de  $K$

Calculer la matrice  $M$  de l'opérateur de bordure.

```

for  $j \in \{0, \dots, m\}$  do
  while  $\exists j' < j, l(j') = l(j)$  do
     $C_j \leftarrow C_j + C_{j'} \bmod 2$ 
  return Paires  $(l(j), j)$ 
```

---

Dans le pire des cas, on a un algorithme en  $\mathcal{O}(m^3)$ .

*Démonstration.* À chaque étape de l'algorithme, la colonne  $C_j$  représente une chaîne de la forme :

$$\partial \left( \sigma^j + \sum_{i < j} \varepsilon_i \sigma^i \right), \varepsilon_i \in \{0, 1\}$$

À la fin de l'algorithme, si  $j$  est tel que  $l(j)$  est défini, alors  $\sigma^{l(j)}$  est un simplexe positif. Donc si à la fin de l'algorithme,  $C_j$  est nulle alors  $\sigma^j$  est positif. Donc, si  $C_j$  n'est pas nulle, alors  $(\sigma^{l(j)}, \sigma^j)$  est une paire de persistence. ■

**Définition 2.7** On représente sur un *diagramme de persistence* les *paires de persistence*  $(\sigma^{l(j)}, \sigma^j)$  par  $(l(j), j)$  ou  $(f(\sigma^{l(j)}), f(\sigma^j))$ . On ajoute au diagramme la diagonale  $\{y = x\}$  et, pour chaque simplexe positif qui n'est pas dans une paire  $\sigma^i$ , le point  $(i, +\infty)$ .

**Définition 2.8** Si  $D_1, D_2$  sont deux diagrammes (potentiellement augmentés pour avoir le même cardinal) :

**La Distance du Goulot** est définie par :

$$d_B^\infty(D_1, D_2) = \inf_{\gamma \in \Gamma} \sup_{p \in D_1} \|p - \gamma(p)\|_\infty$$

**La Distance  $p$ -Wasserstein** est définie, pour  $p \geq 1$  par :

$$W_p(D_1, D_2) = \inf_{\gamma \in \Gamma} \left( \sum_{\rho \in D_1} \|\rho - \gamma(\rho)\|_p^p \right)^{\frac{1}{p}}$$

Dans les deux cas,  $\Gamma$  est l'ensemble des bijections entre  $D_1$  et  $D_2$ .

**R** Ces deux définitions peuvent être vues comme le coût du transport optimal pour la norme infinie et la norme  $p$ .

Ces deux définitions sont par ailleurs équivalentes à un facteur près, ce qui est important pour les théorèmes de stabilité.

**Théorème 2.3** Si  $f, g : X \rightarrow \mathbb{R}$  sont *dociles*, on a :

$$d_B^\infty(D_f, D_g) \leq \|f - g\|_\infty$$

où  $D_\varphi$  est le diagramme de persistance de la filtration associée aux ensembles de sous-niveau de  $\varphi$  sur  $X$ .

On reviendra plus tard sur la notion de docilité.

## 2.3 Calcul de filtrations

### 2.3.1 Complexes de Čech, de Vietoris-Rips, et autres

**Définition 2.9** On considère un recouvrement  $\mathcal{U}$  par des ouverts d'un espace topologique  $X$ . Le *complexe de Čech*  $C(\mathcal{U})$  associé au recouvrement  $\mathcal{U}$  vérifie :

- L'ensemble de sommets de  $C(\mathcal{U})$  est l'ensemble  $\mathcal{U}$ .
- $[U_0, \dots, U_k]$  est un  $k$ -simplexe dans  $C(\mathcal{U})$  si et seulement si  $\cap U_j \neq \emptyset$ .

**Théorème 2.4 — Nerveux (Leray)** Si toutes les intersections entre les ouverts de  $\mathcal{U}$  sont soit vides soit contractibles, alors  $C(\mathcal{U})$  et  $X$  sont homotopiquement équivalents.

Si on se donne plutôt un nuage de point  $V$  dans un espace métrique  $(X, d)$  et un réel  $\alpha$ .

**Définition 2.10** Le *complexe de Čech*  $\check{\text{C}}\text{ech}(V, \alpha)$  est le complexe simplicial filtré indexé par  $\mathbb{R}$  dont l'ensemblde sommets est  $V$  et tel que :

$$\sigma = [p_0, \dots, p_k] \in \check{\text{C}}\text{ech}(V, \alpha) \Leftrightarrow \bigcap_{i=0}^k B(p_i, \alpha) \neq \emptyset$$

**Définition 2.11** Le *complexe de Vietoris-Rips*  $\text{Rips}(V)$  est le complexe simplicial filtré indexé par  $\mathbb{R}$  dont l'ensemble de sommets est  $V$  et est défini par :

$$\sigma = [p_0, \dots, p_k] \in \check{\text{C}}\text{ech}(V, \alpha) \Leftrightarrow \forall i, j \in \{0, \dots, k\}, d(p_i, p_j) \leq \alpha$$

**Proposition 2.4** On a, pour tout  $\alpha > 0$  :

$$\check{\text{Cech}}\left(L, \frac{\alpha}{2}\right) \subseteq \text{Rips}(L, \alpha) \subseteq \check{\text{Cech}}(L, \alpha)$$

**Définition 2.12** Si  $V = \{p_1, \dots, p_n\} \subseteq \mathbb{R}^d$ , on définit la *cellule de Voronoï* associée à  $p_i$  par :

$$\mathcal{V}\nabla(p_i) = \{x \in \mathbb{R}^d \mid \forall j, \|x - p_i\| \leq \|x - p_j\|\}$$

Le *complexe de Delaunay*  $\mathcal{D}(P)$  est le nerf de la couverture faite par les cellules de Voronoï. L'*alpha complexe*  $\mathcal{A}(P, \alpha)$ , pour  $\alpha \geq 0$  est le nerf de la famille :

$$(\text{Vor}(p_i) \cap B(p_i, \sqrt{\alpha}))_{i=1, \dots, n}$$

**Théorème 2.5**  $\mathcal{A}(P, \alpha)$  est homotopie équivalent à  $\bigcup_{i=1}^n B(p_i, \sqrt{\alpha})$ .

### 2.3.2 Stabilité

On va utiliser ci-dessous la distance de Hausdorff :

$$d_H(A, B) = \max \left\{ \sup_{b \in B} d(b, A), \sup_{a \in A} d(a, B) \right\}$$

et la distance de Gromov-Hausdorff :

$$d_{GH}(\mathbb{X}, \mathbb{Y}) = \inf_{\mathbb{Z}, \gamma_1, \gamma_2} (\mathbb{X}, \mathbb{Y})$$

l'infimum étant pris pour  $\mathbb{Z}$  un espace métrique, et  $\gamma_1, \gamma_2$  des immersions isométriques de  $\mathbb{X}, \mathbb{Y}$  dans  $\mathbb{Z}$ .

**Théorème 2.6** Si  $\mathbb{X}$  et  $\mathbb{Y}$  sont des espaces métriques pré-compacts :

$$d_\infty(\text{Rips}(\mathbb{X}), \text{Rips}(\mathbb{Y})) \leq d_{GH}(\mathbb{X}, \mathbb{Y})$$

Ceci est notamment utile lorsqu'on considère la classification de formes non-rigides, puisqu'alors celles-ci sont presque isométriques, mais que calculer leur distance de Gromov-Hausdorff est très coûteux. On va désormais essayer de démontrer ces résultats de stabilité :

**Définition 2.13** Un *module de persistance*  $\mathbb{V}$  est une famille d'espaces vectoriels  $(V_a)_{a \in \mathbb{R}}$  et une famille  $v_a^b : V_a \rightarrow V_b, a \leq b$  qui se compose bien et de sorte que  $v_a^a$  soit l'identité.

- Si  $\mathbb{S}$  est un complexe simplicial filtré, les familles  $V_a = H(\mathbb{S}_a)$  et  $v_a^b : H(\mathbb{S}_a) \rightarrow H(\mathbb{S}_b)$  l'application linéaire induites par l'inclusion  $\mathbb{S}_a \hookrightarrow \mathbb{S}_b$  forment un module de persistance.
- Étant donné un espace métrique  $\mathbb{X}$ ,  $H(\text{Rips}(\mathbb{X}))$  est un module de persistance.
- La filtration par les sous-niveaux de  $f$  induit un module de persistance au niveau de l'homologie.

**R** Il faut voir un module de persistance comme un foncteur de la (petite) catégorie associée au poset d'indexation, vers la catégorie des modules sur un anneau  $A$ . Ici, c'est donc un foncteur de  $\mathbb{R}$  vers  $\text{Vect}_{\mathbb{F}_2}$ , puisqu'on considère notre homologie dans  $\mathbb{F}_2$ .

**Définition 2.14** On dit qu'un module de persistance est dit *q-docile* si pour tout  $a < b$ ,  $v_a^b$  est de rang fini.

Si  $\mathbb{X}$  est pré-compact métrique, alors  $H(\text{Rips}(\mathbb{X}))$  et  $H(\check{\text{Cech}}(\mathbb{X}))$  sont *q-dociles*.

Cette condition apporte de forts théorèmes :

**Théorème 2.7** Les modules de persistence  $q$ -dociles ont des diagrammes de persistence bien définis.

Il faut ici entendre la notion de diagrammes de persistence comme définis par les bases des espaces.

**Définition 2.15** Un *homomorphisme de degré  $\varepsilon$*  entre deux modules de persistence est une collection  $\Phi$  d'application linéaire vérifiant :

$$\begin{array}{ccc} U_a & \xrightarrow{u_a^b} & U_b \\ & \searrow \varphi_a & \swarrow \varphi_b \\ & V_{a+\varepsilon} & \xrightarrow[v_{a+\varepsilon}^{b+\varepsilon}]{} V_{b+\varepsilon} \end{array}$$

Un  $\varepsilon$ -intercalaire entre  $\mathbb{U}$  et  $\mathbb{V}$  est défini par deux homomorphismes de degré  $\varepsilon$   $\Phi : \mathbb{U} \rightarrow \mathbb{V}$  et  $\Psi : \mathbb{V} \rightarrow \mathbb{U}$  de vérifiant :

$$\begin{array}{ccccccc} \dots & \longrightarrow & U_a & \xrightarrow{u_a^{a+2\varepsilon}} & U_{a+2\varepsilon} & \longrightarrow & \dots \\ & \nearrow & \downarrow \varphi_a & \nearrow \psi_{a+\varepsilon} & \downarrow \varphi_{a+2\varepsilon} & \nearrow & \\ \dots & \longrightarrow & V_{a+\varepsilon} & \xrightarrow[v_{a+\varepsilon}^{a+3\varepsilon}]{} & V_{a+3\varepsilon} & & \end{array}$$

**Théorème 2.8** Si  $\mathbb{U}$  et  $\mathbb{V}$  sont  $q$ -dociles et  $\varepsilon$ -intercalés pour un certain  $\varepsilon \geq 0$ , alors :

$$d_\infty(\text{diag}(\mathbb{U}), \text{diag}(\mathbb{V})) \leq \varepsilon$$

On va donc chercher à construire des filtrations qui induisent par leurs groupes d'homologie des modules de persistence  $q$ -dociles, et qui sont  $\varepsilon$ -intercalés quand les espaces/fonctions considérées sont  $O(\varepsilon)$ -proches. Plus particulièrement, on va démontrer la docilité des complexes de Rips et de Čech.

### 2.3.3 Théorèmes de Stabilité

**Définition 2.16** Une *application multivaluée*  $C$  de  $\mathbb{X}$  dans  $\mathbb{Y}$  est une partie de  $\mathbb{X} \times \mathbb{Y}$  qui se projette surjectivement sur  $\mathbb{X}$  par la projection  $\pi_{\mathbb{X}}$  définissant le produit. *L'image*  $C(\sigma)$  de  $\sigma \subseteq \mathbb{X}$  est la projection canonique sur  $\mathbb{Y}$  de la préimage de  $\sigma$  par  $\pi_{\mathbb{X}}$ . La *transposée*  ${}^t C$  de  $C$  est l'image de  $C$  par la symétrie  $\mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{Y} \times \mathbb{X}$ . Si  ${}^t C$  est aussi une application multivaluée, on dit que  $C$  est une *correspondance*.

**Définition 2.17** Si  $(\mathbb{X}, \rho_{\mathbb{X}}$  et  $(\mathbb{Y}, \rho_{\mathbb{Y}})$  sont des espaces métriques compacts, une correspondance  $C$  de  $\mathbb{X}$  dans  $\mathbb{Y}$  est une  $\varepsilon$ -correspondance si :

$$\forall (x, y), (x', y') \in C, |\rho_X(x, x') - \rho_Y(y, y')| \leq \varepsilon$$

**Proposition 2.5** Avec les hypothèses de la définition ci-dessus :

$$d_{GH}(\mathbb{X}, \mathbb{Y}) = \frac{1}{2} \inf \{ \varepsilon \geq 0 \mid \exists C, C \text{ est une } \varepsilon\text{-correspondance} \}$$

**Définition 2.18** Si  $\mathbb{S}, \mathbb{T}$  sont des complexes simpliciaux filtrés avec des ensembles de sommets  $\mathbb{X}$  et  $\mathbb{Y}$  respectivement, une application multivaluée  $C$  de  $\mathbb{X}$  dans  $\mathbb{Y}$  est dite  $\varepsilon$ -simpliciale de  $\mathbb{S}$  dans  $\mathbb{T}$  si pour tout  $a \in \mathbb{R}$  et tout simplexe  $\sigma \in \mathbb{S}_a$ , chaque partie finie de  $C(\sigma)$  est un simplexe de  $\mathbb{T}_{a+\varepsilon}$ .

**Proposition 2.6** Si  $C$  est une correspondance telle que  $C$  et  ${}^t C$  sont  $\varepsilon$ -simpliciales de  $\mathbb{S}$  dans  $\mathbb{T}$ , elle permette de construire un  $\varepsilon$ -intercalaire entre  $H(\mathbb{S})$  et  $H(\mathbb{T})$ .

*Démonstration.* Il suffit pour ça de voir qu'une fonction  $\varepsilon$ -simpliciale définit immédiatement un homomorphisme de degré  $\varepsilon$  sur les modules de persistence définis par les filtrations considérées, par continuation linéaire. ■

**Proposition 2.7** Si  $\mathbb{X}, \mathbb{Y}$  sont des espaces métriques, pour tout  $\varepsilon > 2d_{GH}(\mathbb{X}, \mathbb{Y})$ , alors  $H(\text{Rips}(\mathbb{X}))$  et  $H(\text{Rips}(\mathbb{Y}))$  sont  $\varepsilon$ -intercalés.

*Démonstration.* Si  $C$  est une correspondance de  $\mathbb{X}$  dans  $\mathbb{Y}$  avec une distortion d'au plus  $\varepsilon$ . Si  $\sigma \in \text{Rips}(\mathbb{X}, a)$  alors  $\rho_{\mathbb{X}}(x, x') \leq a$  pour tout  $x, x' \in \sigma$ . Soit  $\tau \subseteq C(\sigma)$  fini. Pour tout  $y, y' \in \tau$ , il existe  $x, x' \in \sigma$  tels que :  $y \in C(x)$  et  $y' \in C(x')$  donc :

$$\rho_{\mathbb{Y}}(y, y') \leq \rho_{\mathbb{X}}(x, x') + \varepsilon \leq a + \varepsilon$$

Ainsi,  $\tau \in \text{Rips}(\mathbb{Y}, a + \varepsilon)$ . De même  ${}^tC$  est  $\varepsilon$ -simpliciale de  $\text{Rips}(\mathbb{Y})$  dans  $\text{Rips}(\mathbb{X})$ . On conclut par la Proposition 2.6. ■

**Proposition 2.8** Si  $\mathbb{X}, \mathbb{Y}$  sont des espaces métriques, pour tout  $\varepsilon \geq 2d_{GH}(\mathbb{X}, \mathbb{Y})$ , alors  $H(\check{\text{Cech}}(\mathbb{X}))$  et  $H(\check{\text{Cech}}(\mathbb{Y}))$  sont  $\varepsilon$ -intercalés.

La preuve est similaire à celle d'avant.

**Théorème 2.9** Soit  $\mathbb{X}$  un espace métrique compact. Les modules de persistence associés à l'homologie des complexes  $\text{Rips}(\mathbb{X})$  et  $\check{\text{Cech}}(\mathbb{X})$  sont  $q$ -dociles.

*Démonstration.* On veut montrer que  $I_a^b : H(\text{Rips}(\mathbb{X}, a)) \rightarrow H(\text{Rips}(\mathbb{X}, b))$  sont de rang finis quand  $a < b$ . On pose  $\varepsilon = (b - a)/2$  et  $F \subseteq \mathbb{X}$  un ensemble fini tel que  $d_H(\mathbb{X}, F) \leq \varepsilon/2$ . Alors :

$$C = \left\{ (x, f) \in X \times F \mid d(x, f) \leq \frac{\varepsilon}{2} \right\}$$

définit une  $\varepsilon$ -correspondance. Utilisant l'application d'intercalage de  $X$  et  $F$ ,  $I_a^b$  se factorise en :

$$H(\text{Rips}(X, a)) \longrightarrow H(\text{Rips}(F, a + \varepsilon)) \rightarrow H(\text{Rips}(X, a + 2\varepsilon)) = H(\text{Rips}(X, b))$$

de dimension finie

**Théorème 2.10** Si  $\mathbb{X}, \mathbb{Y}$  sont des espaces métriques compacts, alors :

$$\begin{aligned} d_\infty(\text{diag}(H(\check{\text{Cech}}(\mathbb{X}))), \text{diag}(H(\check{\text{Cech}}(\mathbb{Y})))) &\leq 2d_{GH}(\mathbb{X}, \mathbb{Y}) \\ d_\infty(\text{diag}(H(\text{Rips}(\mathbb{X}))), \text{diag}(H(\text{Rips}(\mathbb{Y})))) &\leq 2d_{GH}(\mathbb{X}, \mathbb{Y}) \end{aligned}$$

La preuve des deux derniers théorèmes n'utilise pas l'inégalité triangulaire on pourrait donc étendre les résultats précédents à des espaces munies d'une similarité.

Cependant, on a des problèmes avec la dimension de nos espaces : pour tout  $0 < \alpha \leq \beta \in \mathbb{R}$ , il existe un espace métrique compact  $X$  (immersible dans  $\mathbb{R}^4$ ) tel que pour tout  $a \in [\alpha, \beta]$ ,  $H_k(\text{Rips}(X, a))$  est de dimension indénombrable. Toutefois :

- Si  $X$  est compact,  $\dim H_1(\check{\text{Cech}}(X, a)) < +\infty$
- Si  $X$  est géodésique,  $\dim H_1(\text{Rips}(X, a)) < +\infty$  pour  $a > 0$  et  $\text{diag}(H_1(\text{Rips}(X)))$  est contenu dans la ligne  $x = 0$
- Si  $X$  est un espace géodésique  $\delta$ -hyperbolique,  $\text{diag}(H_2(\text{Rips}(X)))$  est contenu dans une bande verticale de largeur  $\mathcal{O}(\delta)$ .

## 2.4 Calculabilité et bruit

Le complexe de Vietoris-Rips et ses filtrations se calculent en  $\mathcal{O}(|\mathbb{X}|^d)$ , ce qui rend le calcul de persistance quasi impossible en pratique. Par ailleurs, les filtrations et la distance de Gromov-Hausdorff sont très sensibles au bruit et aux anomalies.

### 2.4.1 Calcul statistique

On va s'intéresser à un espace métrique  $(\mathbb{M}, \rho)$  et à une mesure de probabilité  $\mu$  à support compact  $X_\mu$  dans  $\mathbb{M}$ . On échantillonne  $m$  points selon  $\mu$ , ce qui nous donne un nuage de point  $\hat{\mathbb{X}}_m$ , et une filtration  $\text{Filt } \hat{\mathbb{X}}_m$ . On a alors :

**Proposition 2.9** Si  $\varepsilon > 0$  :

$$\mathbb{P}\left(d_\infty\left(\text{diag}(\text{Filt}(\mathbb{X}_\mu)), \text{diag}(\text{Filt}(\hat{\mathbb{X}}_m))\right) > \varepsilon\right) \leq \mathbb{P}\left(d_{GH}(\mathbb{X}_\mu, \hat{\mathbb{X}}_m) > \frac{\varepsilon}{2}\right)$$

*Démonstration.* Conséquence directe du Théorème 2.8 de stabilité. ■

On obtient quasi immédiatement des inégalités de déviation :

**Définition 2.19** Pour  $a, b > 0$ , on dit que  $\mu$  vérifie *la supposition  $(a, b)$ -standard* si pour  $x \in \mathbb{X}_\mu$  et  $r > 0$ , on a :

$$\mu(B(x, r)) \geq \min(ar^b, 1)$$

On note  $\mathcal{P}(a, b, \mathbb{M})$  *l'ensemble des distributions de probabilité  $(a, b)$ -standard* sur  $\mathbb{M}$ .

**Théorème 2.11** Si  $\mu$  vérifie la supposition  $(a, b)$ -standard, pour tout  $\varepsilon > 0$  :

$$\mathbb{P}\left(d_\infty\left(\text{diag}(\text{Filt}(\mathbb{X}_\mu)), \text{diag}(\text{Filt}(\hat{\mathbb{X}}_m))\right) > \varepsilon\right) \leq \min\left(\frac{8^b}{a\varepsilon^b} \exp(-ma\varepsilon^b), 1\right)$$

De plus :

$$\mathbb{P}\left(d_\infty\left(\text{diag}(\text{Filt}(\mathbb{X}_\mu)), \text{diag}(\text{Filt}(\hat{\mathbb{X}}_m))\right) \geq C_1 \left(\frac{\log m}{m}\right)^{1/b}\right) \xrightarrow[m \rightarrow \infty]{} 1$$

où  $C_1$  est une constante qui ne dépend que de  $a$  et  $b$ .

*Démonstration.* On commence par majorer  $\mathbb{P}\left(d_{GH}(\mathbb{X}_\mu, \hat{\mathbb{X}}_m) > \frac{\varepsilon}{2}\right)$  puis, on obtient par la supposition  $(a, b)$ -standard une borne supérieure explicite pour la couverture de  $\mathbb{X}_\mu$  par des boules de rayon  $\varepsilon/2$ . On peut alors conclure en prenant l'union des bornes. ■

**Théorème 2.12** On a :

$$\sup_{\mu \in \mathcal{P}(a, b, \mathbb{M})} \mathbb{E}[d_\infty(\text{diag}(\text{Filt}(\mathbb{X}_\mu)), \text{diag}(\text{Filt}(\hat{\mathbb{X}}_m)))] \leq C\left(\frac{\ln m}{m}\right)^{1/b}$$

où  $C$  ne dépend que de  $a$  et  $b$ . Si de plus il y a un point non isolé  $x$  dans  $\mathbb{M}$ , et si  $x_m \in \mathbb{M} \setminus \{x\}$ , telle que  $\rho(x, x_m) \leq (am)^{-1/b}$ , pour tout estimateur  $\hat{\text{diag}}_m$  de  $\text{diag}(\text{Filt}(\mathbb{X}_\mu))$  :

$$\liminf_{m \rightarrow \infty} \rho(x, x_m)^{-1} \sup_{\mu \in \mathcal{P}(a, b, \mathbb{M})} \mathbb{E}[d_\infty(\text{diag}(\text{Filt}(\mathbb{X}_\mu)), \hat{\text{diag}}_m)] \geq C'$$

où  $C'$  est une constante absolue.

### 2.4.2 Paysages de persistence et rééchantillonnage

**Définition 2.20** Si on a un diagramme de persistence  $(b_i, d_i)$ , son *paysage de persistence* est obtenu en y ajoutant à chaque point les deux projections orthogonales sur la diagonale par rapport aux axes, puis en plaçant à l'horizontale sa diagonale. Formellement, c'est l'union pour  $p = (\frac{b+d}{2}, \frac{d-b}{2})$  des graphes :

$$\Lambda_p(t) = \begin{cases} t - b & t \in [b, \frac{b+d}{2}] \\ d - t & t \in [\frac{b+d}{2}, d] \\ 0 & \text{sinon} \end{cases}$$

C'est un encodage de la persistence comme un élément d'un espace fonctionnel, comme fait ci-dessous :

**Définition 2.21** On définit le *k-ème paysage* d'un diagramme  $D$  par :

$$\lambda_D(k, t) = \underset{p \in D}{\text{kmax}} \Lambda_p(t), t \in \mathbb{R}, k \in \mathbb{N}$$

où  $\text{kmax}$  désigne le  $k$ -ème plus grand élément d'un ensemble.

**Proposition 2.10**

- Pour  $t \in \mathbb{R}$  et  $k \in \mathbb{N}$ ,  $0 \leq \lambda_D(k, t) \leq \lambda_D(k+1, t)$

- Pour  $t \in \mathbb{R}$ , et  $k \in \mathbb{N}$ ,  $|\lambda_D(k, t) - \lambda_{D'}(k, t)| \leq d_\infty(D, D')$

Dans la suite, on note  $\mathcal{L}_T$  les paysages dont le support est dans  $[0, T]$ , on prend  $P$  une distribution de probabilité sur  $\mathcal{L}_T$  et  $\lambda_1, \dots, \lambda_n \sim P$  i.i.d. On note  $\mu(t) = \mathbb{E}[\lambda_i(t)]$  le paysage moyen et on l'estime par la moyenne échantillonnée :

$$\bar{\lambda}_n(t) = \frac{1}{n} \sum \lambda_i(t)$$

$\bar{\lambda}_n$  est un estimateur point à point non biaisé de  $\mu$ , qui converge point à point.

**Définition 2.22** Soit  $\mathcal{F}$  la famille des applications d'évaluation  $f_t : \mathcal{L}_T \rightarrow \mathbb{R}$ . Le *processus empirique* indexé par les  $f_t$  est défini par :

$$\mathbb{G}_n(t) = \sqrt{n}(\bar{\lambda}_n(t) - \mu_t) = \sqrt{n}(P_n - P)(f_t)$$

**Théorème 2.13** Soit  $\mathbb{G}$  un pont Brownien avec fonction de covariance :

$$\kappa(s, t) = \int f_t(\lambda) f_s(\lambda) dP(\lambda) - \int f_t(\lambda) dP(\lambda) \int f_s(\lambda) dP(\lambda).$$

On a alors :

$$\mathbb{G}_n \rightarrow \mathbb{G}$$

pour la convergence faible.

Si de plus on note  $\sigma(t)$  l'écart-type de  $\sqrt{n}\bar{\lambda}_n(t)$  :

**Théorème 2.14** Si  $\sigma(t) > c > 0$  sur un intervalle  $I = [t_*, t^*] \subseteq [0, T]$  pour une constante  $c$ , avec  $W = \sup_{t \in I} |\mathbb{G}(f_t)|$  on a :

$$\sup_{z \in \mathbb{R}} \left| \mathbb{P} \left( \sup_{t \in [t_*, t^*]} |\mathbb{G}_n(t)| \leq z \right) - \mathbb{P}(W \leq z) \right| = \mathcal{O} \left( \frac{(\log n)^{7/8}}{n^{1/8}} \right)$$

C'est une forme de théorème central limite uniforme. On a de plus le corollaire suivant :

**Théorème 2.15** Sous les mêmes hypothèses, étant donné un niveau de confiance  $1 - \alpha$ , on peut construire des fonctions de confiance  $l_n(t)$  et  $u_n(t)$  telles que :

$$\mathbb{P}(l_n(t) \leq \mu(t) \leq u_n(t), \forall t \in I) \geq 1 - \alpha - \mathcal{O}\left(\frac{(\log n)^{7/8}}{n^{1/8}}\right)$$

De plus, on a :

$$\sup_t u_n(t) - l_n(t) = \mathcal{O}\left(\sqrt{\frac{1}{n}}\right)$$

Autrement dit, le rééchantillonnage (ou *bootstrap*) permet d'obtenir des intervalles de confiance pour les paysages.

#### 2.4.3 Bruit et méthode d'échantillonnage

On va ici s'intéresser à l'impact de la procédure d'échantillonnage. On rappelle que la définition des distances de  $p$ -Wasserstein peuvent se poser pour n'importe quelles deux mesures de probabilité sur un même espace métrique  $(\mathbb{M}, \rho)$ . On a notamment le théorème suivant :

**Théorème 2.16** Si  $\mu, \nu$  sont des mesures de probabilité sur un même espace métrique  $(\mathbb{M}, \rho)$ , on a :

$$\|\Lambda_{\mu, m} - \Lambda_{\nu, m}\|_\infty \leq m^{\frac{1}{p}} W_p(\mu, \nu)$$

où  $W_p$  dénote la distance de Wasserstein associée à la fonction de coût  $\rho(\cdot, \cdot)^p$ .

Ceci nous assure de l'utilisabilité des méthodes d'échantillonnage, et notamment de la robustesse aux échantillons peu probables et des méthodes sous-échantillonnant.

Avant de démontrer ceci, donnons trois courts lemmes sur les passages des espaces d'échantillonnage aux espaces de paysages :

**Lemme 2.2** Pour toutes mesures de probabilité  $\mu, \nu$  sur  $(\mathbb{M}, \rho)$ , si  $\rho_m$  est une métrique sur  $\mathbb{M}^m$  telle que :

$$\rho_m(X, Y) \leq \left( \sum_{i=1}^m \rho(x_i, y_i)^p \right)^{\frac{1}{p}}$$

alors :

$$W_p(\mu^{\otimes m}, \nu^{\otimes m}) \leq m^{\frac{1}{p}} W_p(\mu, \nu)$$

*Démonstration.* Si  $\Pi$  est un plan de transport entre  $\mu$  et  $\nu$ , alors  $\Pi^{\otimes m}$  est un plan de transport entre  $\mu^{\otimes m}$  et  $\nu^{\otimes m}$  et donc :

$$\begin{aligned} \int_{\mathbb{M}^{2m}} \rho_m(X, Y)^p d\Pi^{\otimes m}(X, Y) &\leq \int_{\mathbb{M}^m \times \mathbb{M}^m} \sum_{i=1}^m \rho(x_i, y_i)^p d\Pi(x_1, y_1) \cdots d\Pi(x_m, y_m) \\ &= m \int_{\mathbb{M} \times \mathbb{M}} \rho(x_1, y_1)^p d\Pi(x_1, y_1), \end{aligned}$$

ce qui conclut la preuve. ■

**Lemme 2.3** En notant  $\varphi^m : \mathbb{M}^m \rightarrow \mathcal{D}$  la fonction qui à  $X$  associe  $\text{diag}(\text{Filt } X)$  dans l'espace des diagrammes de persistance et si  $\Phi_\mu^m$  est le poussé en avant de  $\mu$  par  $\varphi^m$  :

$$W_p(\Phi_\mu^m, \Phi_\nu^m) \leq W_p(\mu^{\otimes m}, \nu^{\otimes m})$$

*Démonstration.* En notant  $\Delta_m(X, Y) = (\psi(\varphi^m(X)), \psi(\varphi^m(Y)))$ , si  $\Pi$  est un plan entre  $\mu^{\otimes m}$  et  $\nu^{\otimes m}$ , alors le plan  $\Delta_m \sharp \Pi$  poussé en avant de  $\Pi$  est un plan entre  $\Phi_\mu^m$  et  $\Phi_\nu^m$  et on a :

$$\begin{aligned} \int_{\mathcal{D}^2} W_\infty(D_X, D_Y)^p d\Delta_m \sharp \Pi(D_X, D_Y) &= \int_{\mathbb{M}^{2m}} W_\infty(\varphi^m(X), \varphi^m(Y))^p d\Pi(X, Y) \\ &\leq \int_{\mathbb{M}^{2m}} d_H(X, Y)^p d\Pi(X, Y) \quad (2.6) \\ &\leq \int_{\mathbb{M}^{2m}} \rho_m(X, Y)^p d\Pi(X, Y) \end{aligned}$$

■

**Lemme 2.4** En notant  $\psi$  l'application de l'espace des diagrammes vers l'espace des paysages munis de la norme infinie et  $\Psi_\mu^m$  le poussé en avant de  $\varphi_\mu^m$  par  $\psi$  :

$$\left\| \mathbb{E}_{\lambda_X \sim \Psi_\mu^m} [\lambda_X] - \mathbb{E}_{\lambda_Y \sim \Psi_\nu^m} [\lambda_Y] \right\|_\infty \leq W_{\infty, p}(\Phi_\mu^m, \Phi_\nu^m)$$

où  $W_{\infty, p}$  fait appel à la  $p$ -ème puissance de la distance infinie de Wasserstein pour les diagrammes sous-jacents.

*Démonstration.* Si  $\Pi$  est un plan entre  $\Phi_\mu^m$  et  $\Phi_\nu^m$ , pour tout  $t \in \mathbb{R}$  on a :

$$\begin{aligned} \left| \mathbb{E}_{\lambda_X \sim \Psi_\mu^m} [\lambda_X](t) - \mathbb{E}_{\lambda_Y \sim \Psi_\nu^m} [\lambda_Y](t) \right|^p &= |\mathbb{E}[\lambda_X(t) - \lambda_Y(t)]|^p \\ &\leq \mathbb{E}[|\lambda_X(t) - \lambda_Y(t)|^p] \quad (\text{Jensen}) \\ &\leq \mathbb{E}[W_\infty(D_X, D_Y)^p] \quad (2.10) \\ &= \int_{\mathcal{D} \times \mathcal{D}} W_\infty(D_X, D_Y)^p d\Pi(D_X, D_Y) \end{aligned}$$

■

*Preuve du Théorème 2.16.* Par le Lemme 2.2 ci-dessus

$$W_p(\mu^{\otimes m}, \nu^{\otimes m}) \leq m^{\frac{1}{p}} W_p(\mu, \nu)$$

Par le Lemme 2.3, en notant  $P_\pi$  le diagramme de persistence associé à  $\pi$  :

$$W_p(P_\mu, P_\nu) \leq W_p(\mu^{\otimes m}, \nu^{\otimes m})$$

Enfin, par le Lemme 2.4 :

$$\|\Lambda_{\mu, m} - \Lambda_{\nu, m}\|_\infty \leq W_p(P_\mu, P_\nu)$$

■

## 2.5 Persistence et apprentissage automatique

### 2.5.1 Représentation de persistence

Puisque l'espace des diagrammes de persistence n'est pas linéaire, les algorithmes de ML classique ne fonctionnent pas bien. La bibliothèque Python et C++ *Gudhi* propose une large zoologie de représentations pour la persistence, comme mesures discrètes, espaces métriques finis, racines de polynômes ou collections de fonctions 1D.

Par exemple, on peut représenter un diagramme en le plongeant dans  $\mathbb{R}^2$  et l'espace des mesures par  $D = \sum \delta_{p_i}$ . Si on se donne un noyau  $K : \mathbb{R}^2 \rightarrow \mathbb{R}$  et  $H$  une matrice de bande-passante (forme quadratique), en définissant

$K_H(u) = |H|^{-1/2} K(H^{-1/2} \cdot u)$ , on obtient alors, étant donné une fonction de poids  $w$ , *la surface de persistance* de  $D$  par :

$$\forall u \in \mathbb{R}^2, \rho(D)(u) = D(wK_H(u - \cdot))$$

La question se pose alors de savoir comment choisir une représentation adaptée à un réseau de neurones. Une réponse partielle peut être trouvée en regardant l'architecture à ensembles profonds : on se donne  $n$  points dans  $\mathbb{R}^d$  et on construit un réseau dont les niveaux sont invariants par permutation ( $f \circ \sigma = f$ )

**Théorème 2.17 — Universalité** Une fonction  $f$  est invariante par permutation si et seulement si, pour tout  $X$  inclus dans un ensemble dénombrable  $f(X) = \rho(\sum_i \varphi(x_i))$  pour certaines fonctions  $\rho$  et  $\varphi$ .

Les réseaux à niveaux invariants par permutation permettent de généraliser plusieurs approches générales en TDA, sous la forme :

$$\text{PersLay}(\text{diag}) = \rho(\text{op}\{w(p), \varphi(p)\}_{p \in \text{diag}})$$

où  $\text{op}$  est invariante par permutation,  $w$  est une fonction de poids, et  $\varphi$  est une transformation permettant de se ramener à un ensemble dénombrable.

On peut par exemple retrouver la surface de persistance en se donnant  $t_1, \dots, t_q \in \mathbb{R}^2$  puis en posant :

- $w(p) = w_t((x, y))$  ;
- $\varphi_\Gamma : p \mapsto (\Gamma_p(t_i))_i$  avec  $\Gamma_p$  la gaussienne centrée en  $p$  d'écart-type fixé  $\sigma$  ;
- $\text{op} = \sum$ .

Pour les paysages, on prend  $w(p) = 1$ ,  $\text{op} = \text{top}-k$  et  $\varphi_\Lambda$  l'évaluation de  $\Lambda_p$  en  $q$  paramètres  $t_1, \dots, t_q$ .

### 2.5.2 Différentiabilité de la persistance

Nombre de méthodes permettent de minimiser une fonction sur l'ensemble des diagrammes, mais la plupart sont restreintes à un type spécifique de filtration ou de fonction à minimiser.

## 3 Fonctions de Morse

## 4 Inférence Topologique

## 5 Théorie de Morse Discrète

## 6 Noyaux et Statistiques