

Topological data analysis and persistent homology

Developing an outlier detector based on persistent homology

Arnau Mas

Supervised by
Dr Albert Ruiz

Universitat Autònoma de Barcelona
September 3rd 2020



Universitat Autònoma
de Barcelona

- ① The problem
- ② An overview of homology
- ③ The detector
- ④ The results

- Has its origins in the field of oncology: tumour imaging

- Has its origins in the field of oncology: tumour imaging
- Extract many **radiomic features** from scans

- Has its origins in the field of oncology: tumour imaging
- Extract many **radiomic features** from scans
- Look for correlations between radiomic features and treatment response

Benefits

- Non-invasive
- More systematic
- Quantitative

Benefits

- Non-invasive
- More systematic
- Quantitative

Drawbacks

- May be hard to reproduce
- High dimensional feature spaces
- Limited number of cases



- TOPiomics is an ongoing collaboration between the Computer Vision Center (CVC) and the Vall d'Hebron Institute of Oncology (VHIO).



- TOPiomics is an ongoing collaboration between the Computer Vision Center (CVC) and the Vall d'Hebron Institute of Oncology (VHIO).
- Main aim: early detection of patients with strange radiomic signatures, *outliers*, for whom standard treatments could fail.



- TOPiomics is an ongoing collaboration between the Computer Vision Center (CVC) and the Vall d'Hebron Institute of Oncology (VHIO).
- Main aim: early detection of patients with strange radiomic signatures, *outliers*, for whom standard treatments could fail.
- Use *topological data analysis*, hence TOPiomics (topological radiomics).

- Progress thus far: use the *mutual k-nearest neighbours* (MkNN) graph to detect outliers.

- Progress thus far: use the *mutual k-nearest neighbours* (MkNN) graph to detect outliers.
- Problem: there is a parameter that has to be fixed. It has to do with the scale at which we look at the data.

- Progress thus far: use the *mutual k-nearest neighbours* (MkNN) graph to detect outliers.
- Problem: there is a parameter that has to be fixed. It has to do with the scale at which we look at the data.
- Idea: look at *every* parameter value.

- Progress thus far: use the *mutual k-nearest neighbours* (MkNN) graph to detect outliers.
- Problem: there is a parameter that has to be fixed. It has to do with the scale at which we look at the data.
- Idea: look at *every* parameter value.

This is persistence

- ① The problem
- ② An overview of homology
- ③ The detector
- ④ The results

Algebraic topology studies topological spaces by computing algebraic invariants, namely with functors

$$\text{Top} \rightarrow \text{Grp}$$

$$\text{Top} \rightarrow \text{Ab}$$

$$\text{Top} \rightarrow \text{Ring}$$

$$\text{Top} \rightarrow \text{Vect}_K$$

$$\vdots$$

Algebraic topology studies topological spaces by computing algebraic invariants, namely with functors

$$\text{Top} \rightarrow \text{Grp}$$

$$\text{Top} \rightarrow \text{Ab}$$

$$\text{Top} \rightarrow \text{Ring}$$

$$\text{Top} \rightarrow \text{Vect}_K$$

$$\vdots$$

Homology groups are one of these functors,

$$H_n: \text{Top} \rightarrow \text{Ab}.$$

Algebraic topology studies topological spaces by computing algebraic invariants, namely with functors

$$\text{Top} \rightarrow \text{Grp}$$

$$\text{Top} \rightarrow \text{Ab}$$

$$\text{Top} \rightarrow \text{Ring}$$

$$\text{Top} \rightarrow \text{Vect}_K$$

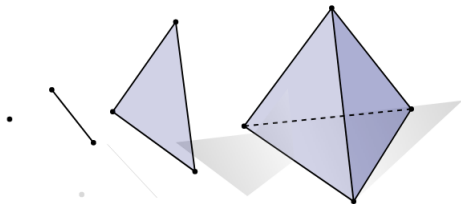
$$\vdots$$

Homology groups are one of these functors,

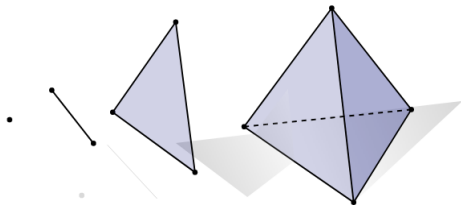
$$H_n: \text{Top} \rightarrow \text{Ab}.$$

They come in many flavours, the simplest one is *simplicial homology*, which works for *simplicial complexes*.

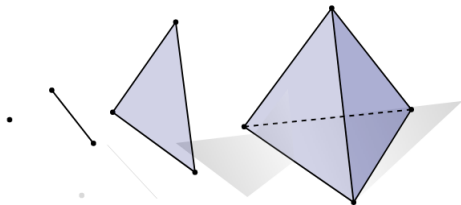
- *Simplex*. The convex hull of a set of points.



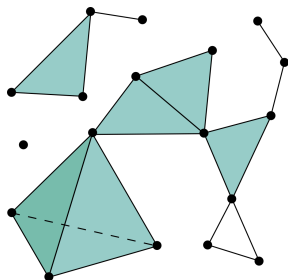
- *Simplex*. The convex hull of a set of (geometrically independent) points.



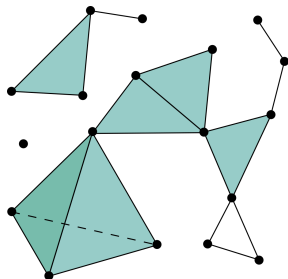
- *Simplex*. The convex hull of an (ordered) set of (geometrically independent) points.



- *Simplex*. The convex hull of an (ordered) set of (geometrically independent) points.
- *Simplicial complex, K* . Topological space assembled out of simplices glued along their faces.



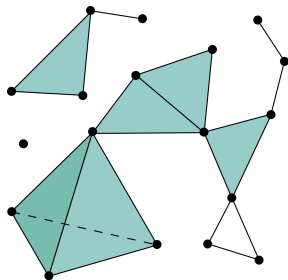
- *Simplex*. The convex hull of an (ordered) set of (geometrically independent) points.
- *Simplicial complex, K* . Topological space assembled out of simplices glued along their faces.
- *Chain groups, $C_n(K)$* . Free abelian groups generated by the n -simplices of a complex.



Remark

In the combinatorial picture, the geometrical requirements are dropped and one speaks of *abstract* simplices, *abstract* simplicial complexes, etc.

- *Simplex*. The convex hull of an (ordered) set of (geometrically independent) points.
- *Simplicial complex, K* .
Topological space assembled out of simplices glued along their faces.
- *Chain groups, $C_n(K)$* . Free abelian groups generated by the n -simplices of a complex.



Remark

Free abelian groups are equivalent to \mathbb{Z} -modules. Alternatively, one can also consider the vector spaces generated by the chains, for instance over \mathbb{F}_2 .

Orientation

The simplex generated by p_0, \dots, p_n is written $[p_0, \dots, p_n]$. The ordering determines an orientation, and we require that flipping orientation implies a change of sign in $C_n(K)$, e.g.

$$[p_0, p_1, \dots, p_n] = -[p_1, p_0, \dots, p_n].$$

Orientation

The simplex generated by p_0, \dots, p_n is written $[p_0, \dots, p_n]$. The ordering determines an orientation, and we require that flipping orientation implies a change of sign in $C_n(K)$, e.g.

$$[p_0, p_1, \dots, p_n] = -[p_1, p_0, \dots, p_n].$$

The chain groups are not really free!

Define the boundary morphism $\partial_n: C_n(K) \rightarrow C_{n-1}(K)$ by

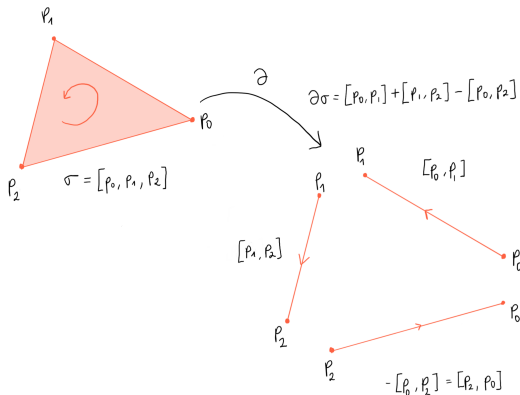
$$\partial_n[p_0, \dots, p_n] := \sum_{k=0}^n (-1)^k [p_0, \dots, \hat{p}_k, \dots, p_n]$$

and extending to chains.

Define the boundary morphism $\partial_n: C_n(K) \rightarrow C_{n-1}(K)$ by

$$\partial_n[p_0, \dots, p_n] := \sum_{k=0}^n (-1)^k [p_0, \dots, \hat{p}_k, \dots, p_n]$$

and extending to chains.



- *Cycles*: chains with no boundary,

$$Z_n(K) := \ker \partial_n \subseteq C_n(K)$$

- *Cycles*: chains with no boundary,

$$Z_n(K) := \ker \partial_n \subseteq C_n(K)$$

- *Boundaries*: chains which are the boundary of a higher-dimensional chain,

$$B_n(K) := \operatorname{im} \partial_{n+1} \subseteq C_n(K)$$

- *Cycles*: chains with no boundary,

$$Z_n(K) := \ker \partial_n \subseteq C_n(K)$$

- *Boundaries*: chains which are the boundary of a higher-dimensional chain,

$$B_n(K) := \operatorname{im} \partial_{n+1} \subseteq C_n(K)$$

Fact

$$\partial_n \circ \partial_{n+1} = 0$$

so boundaries are always cycles.

Homology groups are the quotients

$$H_n(K) := Z_n(K)/B_n(K).$$

Homology groups are the quotients

$$H_n(K) := Z_n(K)/B_n(K).$$

The key insight: Homology groups measure **cycles which are not boundaries**.

Homology groups are the quotients

$$H_n(K) := Z_n(K)/B_n(K).$$

The key insight: Homology groups measure **cycles which are not boundaries**. They tell us about the number of n -dimensional holes in our space.

Homology groups are the quotients

$$H_n(K) := Z_n(K)/B_n(K).$$

The key insight: Homology groups measure **cycles which are not boundaries**. They tell us about the number of n -dimensional holes in our space.

Homology is invariant under homotopy, weaker than invariance under homeomorphism.

Instead of a single simplicial complex, consider a *filtration*,

$$K^0 \subseteq K^1 \subseteq \dots \subseteq K^N.$$

Instead of a single simplicial complex, consider a *filtration*,

$$K^0 \subseteq K^1 \subseteq \dots \subseteq K^N.$$

Then look at this evolution at the level of homology

$$H_{\mathbf{d}}^0(K) \hookrightarrow H_{\mathbf{d}}^1(K) \hookrightarrow \dots \hookrightarrow H_{\mathbf{d}}^N(K)$$

Instead of a single simplicial complex, consider a *filtration*,

$$K^0 \subseteq K^1 \subseteq \dots \subseteq K^N.$$

Then look at this evolution at the level of homology

$$H_d^0(K) \hookrightarrow H_d^1(K) \hookrightarrow \dots \hookrightarrow H_d^N(K)$$

The steps of the filtration can be thought of as time

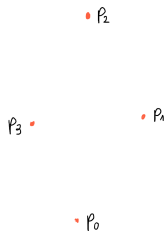
Because of functoriality, there are maps between the different steps of homology,

$$f_i^j: H_d^i(K) \rightarrow H_d^j(K).$$

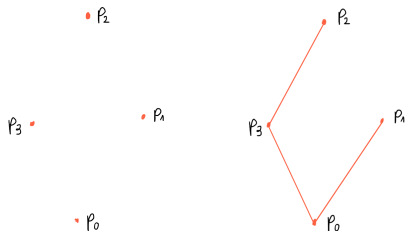
Because of functoriality, there are maps between the different steps of homology,

$$f_i^j : H_d^i(K) \rightarrow H_d^j(K).$$

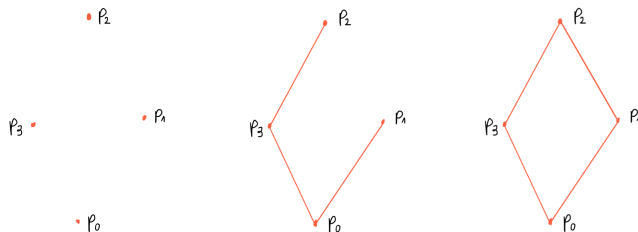
At each step, classes can be born (not in the image of f_i^{i+1}), die (in the kernel of f_i^{i+1}) or merge (have the same image through f_i^{i+1}).



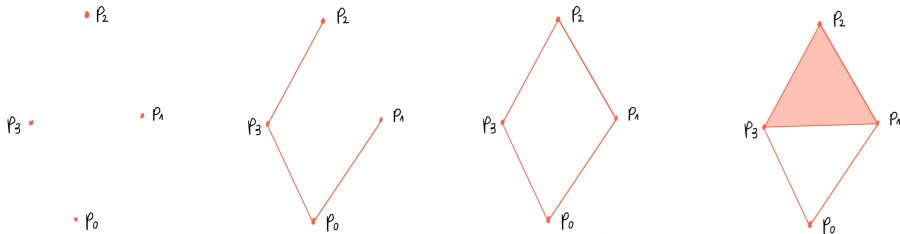
$[p_0], [p_1], [p_2], [p_3]$ are born.



$[p_0], [p_1], [p_2], [p_3]$ merge.



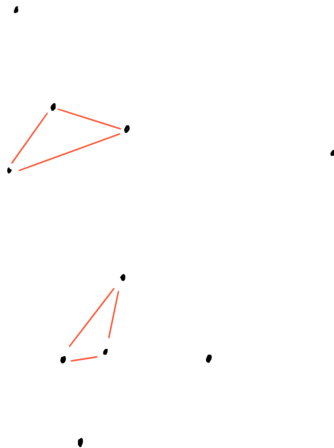
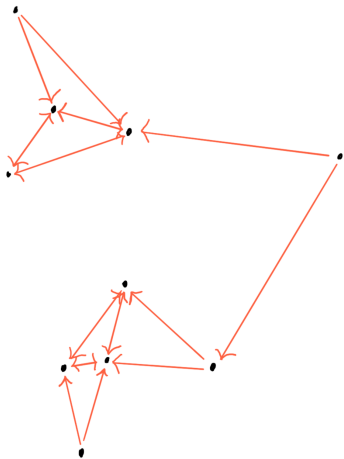
The cycle $[p_0, p_1] + [p_1, p_2] + [p_2, p_3] + [p_3, p_4]$ is born.



The cycle does not die.

- ① The problem
- ② An overview of homology
- ③ The detector
- ④ The results

It is good for detecting clusters in a data set.



The algorithm has four main parts:

The algorithm has four main parts:

- The class `Cloud` stores a point cloud and computes its $MkNN$ graph.

The algorithm has four main parts:

- The class `Cloud` stores a point cloud and computes its $MkNN$ graph.
- The class `Filtration` builds up the filtration.

The algorithm has four main parts:

- The class `Cloud` stores a point cloud and computes its $MkNN$ graph.
- The class `Filtration` builds up the filtration.
- The class `Homology` computes the homology.

The algorithm has four main parts:

- The class `Cloud` stores a point cloud and computes its $MkNN$ graph.
- The class `Filtration` builds up the filtration.
- The class `Homology` computes the homology.
- The persistence data of each class is gathered and plotted to detect the outliers.

The algorithm has four main parts:

- The class `Cloud` stores a point cloud and computes its $MkNN$ graph.
- The class `Filtration` builds up the filtration.
- The class `Homology` computes the homology.
- The persistence data of each class is gathered and plotted to detect the outliers.

All of the Python code is available at
github.com/arnaumas/mknn-homology.

- Store the cloud of N points in \mathbb{R}^d as a NumPy array of shape (N, d) .

- Store the cloud of N points in \mathbb{R}^d as a NumPy array of shape (N, d) .
- Construct the adjacency matrix, A , of the k -nearest neighbours graph of the cloud.

- Store the cloud of N points in \mathbb{R}^d as a NumPy array of shape (N, d) .
- Construct the adjacency matrix, A , of the k -nearest neighbours graph of the cloud.
- The adjacency matrix of the the M k NN graph is

$$M = A \odot A^T$$

where \odot is the entrywise product.

- The package NetworkX can compute the cliques of a graph. These correspond to the simplices of the *graph complex* associated to the $MkNN$ graph of the cloud.

- The package NetworkX can compute the cliques of a graph. These correspond to the simplices of the *graph complex* associated to the $MkNN$ graph of the cloud.
- The simplices are ordered by order of appearance and dimension, which ensures any simplex is always preceded by its faces.

- The classes `Simplex` and `Chain` model elements of the chain groups and underpin the following computations.

- The classes `Simplex` and `Chain` model elements of the chain groups and underpin the following computations.
- The chain group is taken over \mathbb{F}_2 , which corresponds to disregarding orientation and simplifies computation.

- The classes `Simplex` and `Chain` model elements of the chain groups and underpin the following computations.
- The chain group is taken over \mathbb{F}_2 , which corresponds to disregarding orientation and simplifies computation.
- There is a Python dictionary which keeps track of the homology class of every clique. This is a memoised version of the projection $\pi: C_n(K) \rightarrow C_n(K)/B_n(K)$.

- The classes `Simplex` and `Chain` model elements of the chain groups and underpin the following computations.
- The chain group is taken over \mathbb{F}_2 , which corresponds to disregarding orientation and simplifies computation.
- There is a Python dictionary which keeps track of the homology class of every clique. This is a memoised version of the projection $\pi: C_n(K) \rightarrow C_n(K)/B_n(K)$.

- Iterating through the filtration, for every simplex σ , the homology class of the sum of its faces, $w_0 + \cdots + w_n$.

- Iterating through the filtration, for every simplex σ , the homology class of the sum of its faces, $w_0 + \cdots + w_n$.
- If it is zero, then before σ appears, there is chain c such that

$$w_0 + \cdots + w_n = \partial c.$$

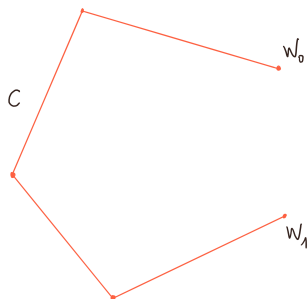
- Iterating through the filtration, for every simplex σ , the homology class of the sum of its faces, $w_0 + \cdots + w_n$.
- If it is zero, then before σ appears, there is chain c such that

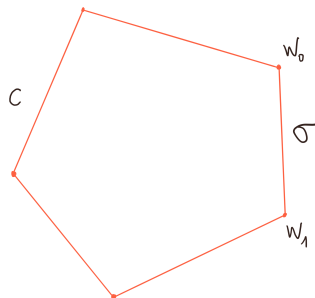
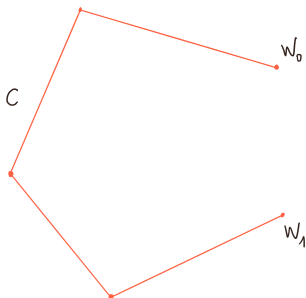
$$w_0 + \cdots + w_n = \partial c.$$

Once σ is added,

$$w_0 + \cdots + w_n = \partial c = \partial \sigma$$

so $\partial(\sigma + c) = 0$, which means $\sigma + c$ is a new cycle.

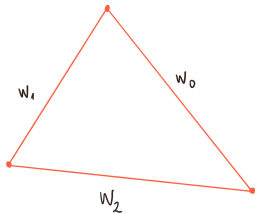




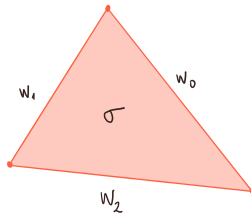
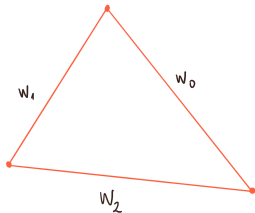
- If it is not zero, then, before σ appears, its faces formed a cycle not homologous to zero. Therefore its homology class was a nontrivial element of the corresponding homology group.

- If it is not zero, then, before σ appears, its faces formed a cycle not homologous to zero. Therefore its homology class was a nontrivial element of the corresponding homology group.
- Once σ appears, it is covered and so it dies.

- If it is not zero, then, before σ appears, its faces formed a cycle not homologous to zero. Therefore its homology class was a nontrivial element of the corresponding homology group.
- Once σ appears, it is covered and so it dies.



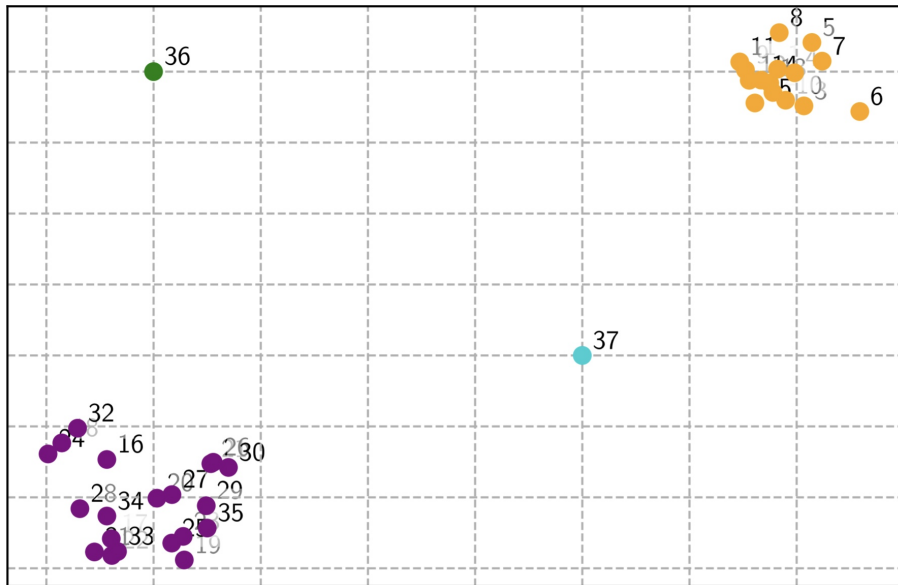
- If it is not zero, then, before σ appears, its faces formed a cycle not homologous to zero. Therefore its homology class was a nontrivial element of the corresponding homology group.
- Once σ appears, it is covered and so it dies.

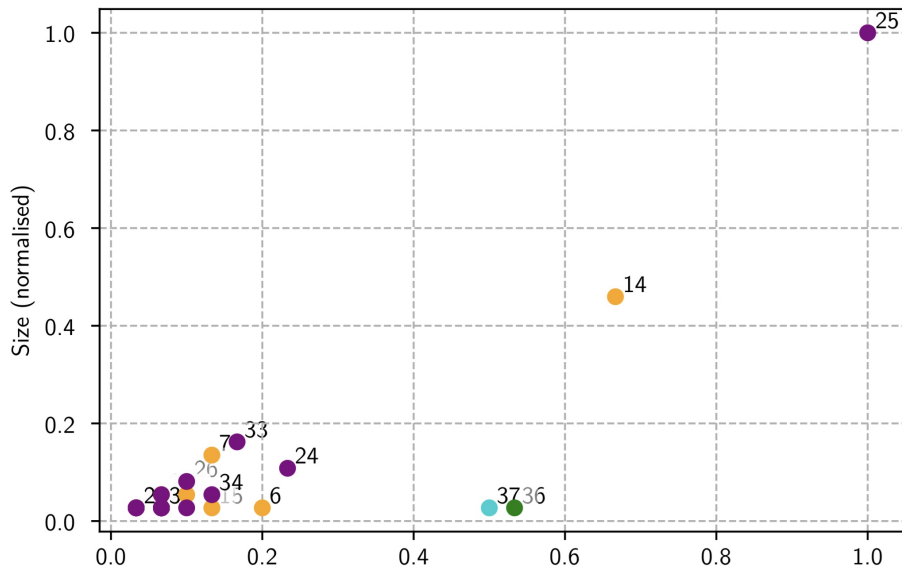


The lifetime and size at death of every class in the zeroth homology group is stored and plotted. Outliers will be those classes with few members and long lifetimes, since they took a long time to merge into larger clusters.

The lifetime and size at death of every class in the zeroth homology group is stored and plotted. Outliers will be those classes with few members and long lifetimes, since they took a long time to merge into larger clusters. Visually, outliers will be those points in the lower right corner of the plot.

- ① The problem
- ② An overview of homology
- ③ The detector
- ④ The results

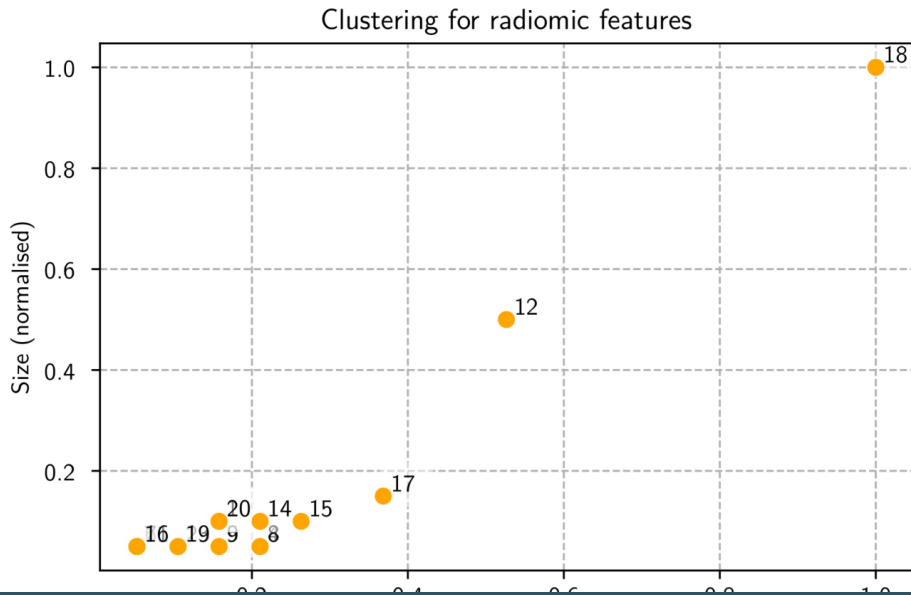


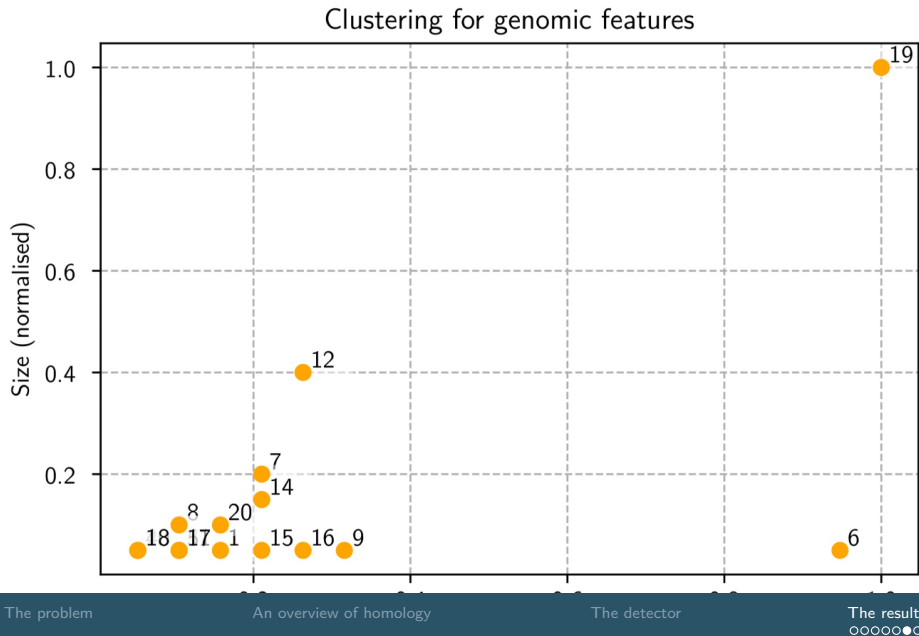


- The detector was run with a data set consisting of 26 radiomic features measured from 20 different patients. These are the result of a previous study carried out at VHIO.

- The detector was run with a data set consisting of 26 radiomic features measured from 20 different patients. These are the result of a previous study carried out at VHIO.
- There are also genomic features corresponding to the same 20 patients.

- The detector was run with a data set consisting of 26 radiomic features measured from 20 different patients. These are the result of a previous study carried out at VHIO.
- There are also genomic features corresponding to the same 20 patients.
- The detector was run on both sets to compare the clustering structure.





- Detector appears to properly detect outliers, and gives a quantitative measure of “outlierness”.

- Detector appears to properly detect outliers, and gives a quantitative measure of “outlierness”.
- The results from the real dataset are awaiting an assesment from VHIO on their medical significance.

- Detector appears to properly detect outliers, and gives a quantitative measure of “outlierness”.
- The results from the real dataset are awaiting an assesment from VHIO on their medical significance.
- Further considerations: efficiency, edgecase testing, higher dimensional homology.

Topological data analysis and persistent homology

Developing an outlier detector based on persistent homology

Arnau Mas

Supervised by
Dr Albert Ruiz

Universitat Autònoma de Barcelona
September 3rd 2020



Universitat Autònoma
de Barcelona