

Topological Data Analysis of Images

October 27, 2020

1. Persistent Homology

A **persistence complex** is a sequence of chain complexes along with chain maps $x_i : C_*^i \longrightarrow C_*^{i+1}$.

For $i < j$ the $(i, j)^{th}$ **persistent homology** of C denoted by $H_*^{i \rightarrow j}(C)$ is the image of $i : H_*(C^i) \longrightarrow H_*(C^j)$.

For a finite persistence module over a field, we can use the structure theorem over PID to interpret the homology module $H_*(C; F)$. The free part is in correspondence with homology generators that appear at a specific value and persist forever, while the torsion part is in correspondence with homology parameters that appear at a particular value and disappear at a higher value.

Thus, the parameter intervals arising from basis for $H_*(C; F)$ in structure theorem can be represented as horizontal line segments ordered arbitrarily. This is called a **persistence barcode**

2. Vectorisation of Persistence Barcodes

Let the barcode be represented by $D = \{(b_j, d_j)\}_{j \in I}$, where I is the set of all bars.

1. Persistent Entropy: Let $l_i = d_i - b_i$, denote the length of the bars in D and $L = \sum l_i$ denote the total length. The persistent entropy of the barcode is the Shannon entropy of the lengths of the bars.

$$PE(D) = \frac{1}{L} \sum l_i \log\left(\frac{l_i}{L}\right) \quad (1)$$

2. Betti Curve: The Betti curve is a real valued function defined on the set of parameter values. At each point, its value is the number of bars that contain this point. The L^p norm of these curves are considered.

2. Vectorisation of Persistence Barcodes (contd.)

3. Persistent Landscapes: A function f , is associated with each barcode in D .

$$f_{(b_i, d_i)}(x) = \begin{cases} 0 & x \notin (b_i, d_i) \\ x - b & x \in (b, \frac{b+d}{2}) \\ -x + d & x \in (\frac{b+d}{2}, d) \end{cases} \quad (2)$$

.

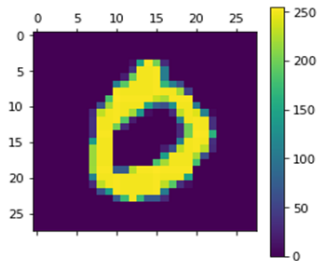
The k -th landscape function is the pointwise k -th maximum of the functions $\{f_{(b_i, d_i)}\}$. The L^p norms of these functions are considered for vectorisation.

4. Wasserstein Amplitude: This is defined as the Wasserstein distance of the given persistence barcode to the empty barcode.

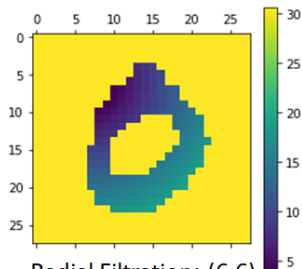
3. Analysis of Images

- ▶ Image classification using topological data analysis involves associating each image with a persistence complex and generating a feature vector corresponding to this. These feature vectors can subsequently be fed to machine learning algorithms
- ▶ Grayscale images can be viewed as a real valued function over a rectangular grid. This structure lends itself for the construction of a sequence of cubical complexes corresponding to the level sets determined by the grayscale filtration function.
- ▶ On binarisation at a suitable threshold, various other filtrations can be defined on these images based on the distribution of the 0 and 1 pixels on the rectangular grid.

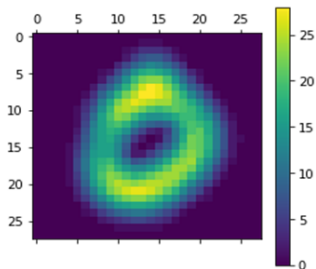
4. Examples of Different Filtrations: MNIST Dataset



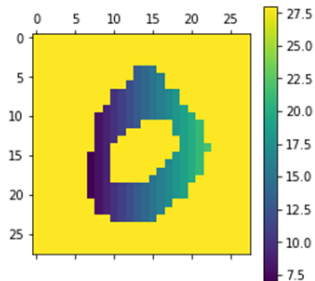
Grayscale Filtration



Radial Filtration: (6,6)



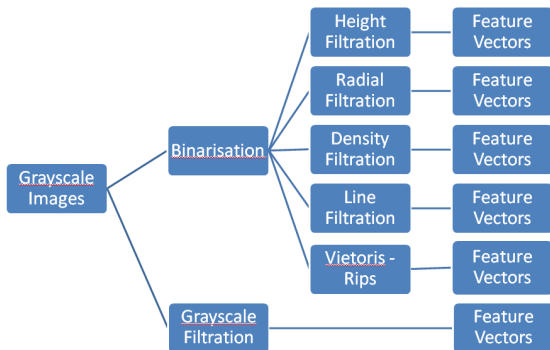
Density Filtration : 4



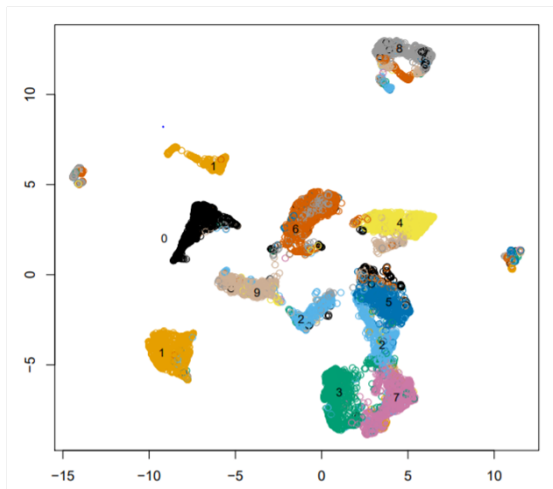
Height Filtration: (1,0)

5. Topological Pipeline

Feature vectors are generated by vectorisation of the persistent barcodes obtained from the persistent homology of the level sets of the filtrations.



6. MNIST Classification



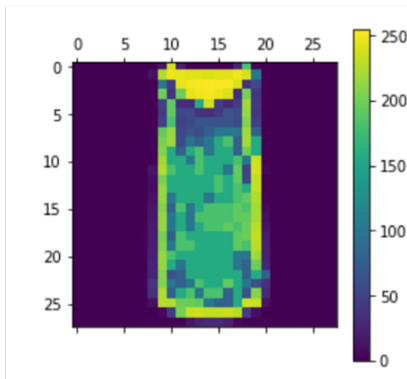
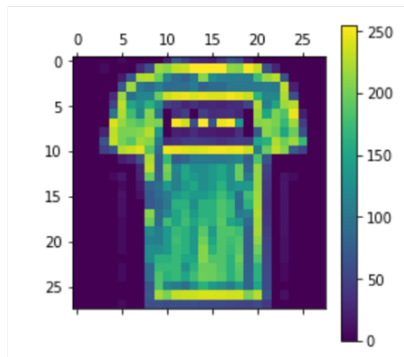
UMAP plot of dimension 2 of the 52 feature vectors generated from the topological pipeline by considering persistent entropy vectorisation.

6. MNIST Classification (contd.)

1. k-NN Classification: Performed with $k = 5$ on 2 dimensional data obtained by using UMAP on feature vectors.
Accuracy : **90.82**
2. Random Forest Classifier: Number of trees = 1000

S.No	Binarisation	Dimension	Filtrations	Vectorisation	Accuracy
1	0.2	50	Height, Radial, Density, Line	Persistent Landscape	94.92
2	0.4	52	Height, Radial, Density, Line, V-R	Entropy	96.15
3	0.3	52	Height, Radial, Density, Line, V-R	Entropy	96.21
4	0.2	52	Height, Radial, Density, Line, V-R	Entropy	96.48
5	0.2	202	Height, Radial, Density, Line, V-R	All Vectorisations	97.16

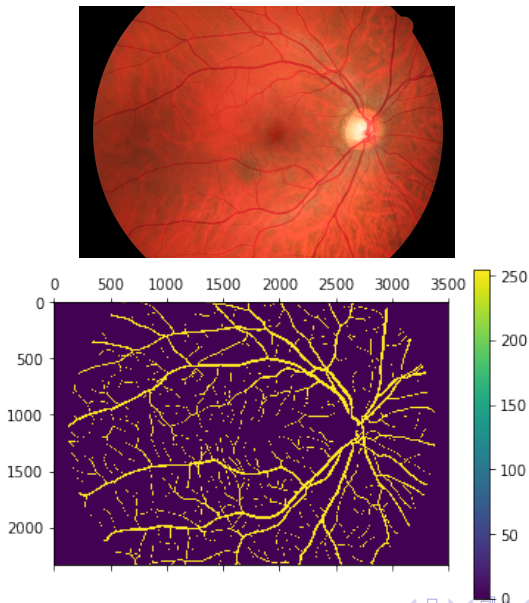
7. Fashion MNIST and Other Datasets



Using 200 feature vectors generated by considering all four vectorisations and height, radial, density and line filtrations of the topological pipeline, an accuracy of **82.85** was obtained on the Fashion MNIST dataset.

7. Fashion MNIST and Other Datasets (contd.)

Retina Fundus Images:



8. References

- [1] R. Ghrist, Barcodes: The persistent topology of data
- [2]. A. Garin and G. Tauzin,"A Topological "Reading" Lesson: Classification of MNIST using TDA",
<https://arxiv.org/pdf/1910.08345.pdf>