

Computational topology: Image classification using Vietoris-Rips complex

May 13, 2016

Neža Mramor Kosta, Gregor Jerše

Andrej Dolenc, Peter Us, Rok Ivanšek

Contents

Project description	3
Obtaining and preprocessing data	3
The classification model	4
Relation to single linkage clustering algorithm	5
Computational complexity	6
Results	6
Cup, paper, pen image set	6
General classification	8
Summary	9
Workload	9



Figure 1: Sample images from used dataset.

Project description

As the title implies, the idea behind the project is to use the Vietoris-Rips (in future abbreviated VR-cx) complex for image classification. We took pictures of three generic objects and preprocessed them so we obtained vector representations of pictures. The vector representations of images can then be looked upon as points in some n -dimensional space \mathbb{R}^n . Over this points we build the VR-cx for two different parameters. We use the complexes to build a classification model for the images and classify them. We also perform several interesting tests on the generated VR-cx.

Obtaining and preprocessing data

For our dataset we used grayscale photographs of 3 different objects: a paper cup, crumpled up paper and a pen. Each of the objects was photographed 10 times from slightly different angles and under different illumination. The sizes of the images were 1080 pixels in width and 810 pixels in height. A selection of these photographs can be seen in figure 1.

We can think of images as matrices of size $(n \times m)$ where each entry represents the intensity of the corresponding pixel (meaning in our case we had 30 matrices of size 1080×810). With images in grayscale color space, each of the values is an integer between 0 (black) and 255 (white). In order to prepare the dataset for evaluation with our classification model, we flattened image's matrix into one long vector of size nm (first row of the matrix being first n elements, second row the next n and so on m times). We then stacked all individual vectors vertically to form a large matrix with 30 rows, one for each individual image, and nm columns, one for each individual pixel.

Because we wanted to help our model as much as possible via preprocessing, we scaled/transformed the values of each pixel either using a standard scaler, or by subtracting mean and scaling values to the $[-1, 1]$ interval. The main motivation behind this was to remove varying illumination from our dataset.

With this however we are still left with a matrix of very high dimensions, so in order to reduce processing time, decrease redundancy and amount of noise, we apply principal component analysis to it. With this we can to reduce the number of columns from $n \times m$ to but a few which makes the dataset much more manageable.

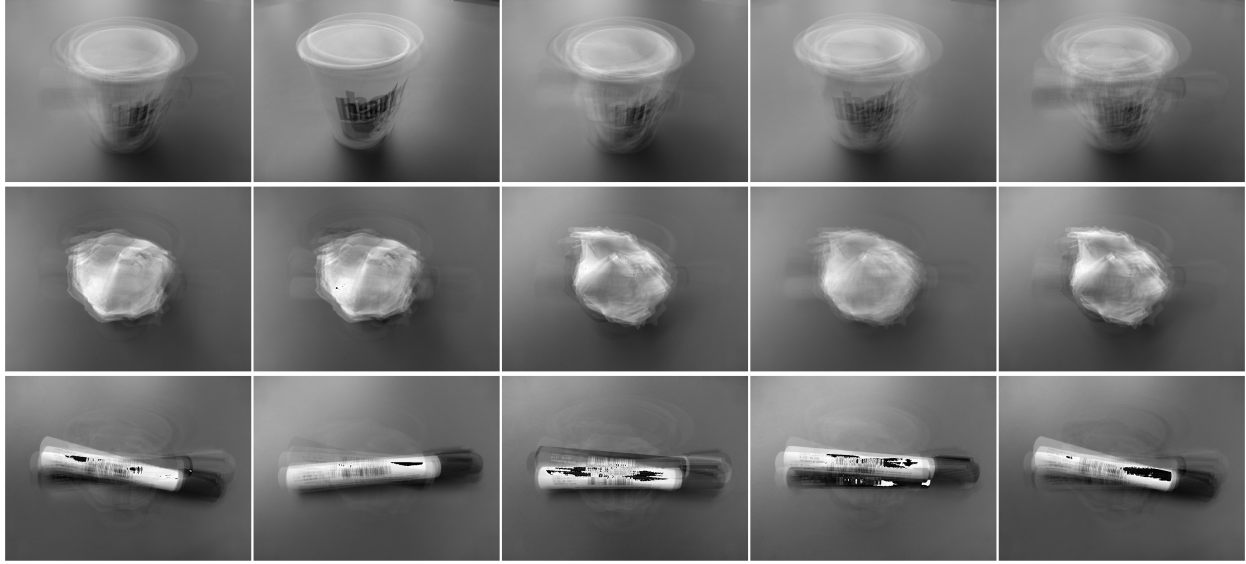


Figure 2: Images from figure 1 after applying preprocessing and inverse preprocessing steps to them. **TODO:** No reference anywhere in the text

The classification model

Definition 1 (Vietoris-Rips complex). *Let X be a set of m -dimensional points $X \in \mathbb{R}^m$ and let d be a metric. Pick a parameter $r > 0$. Construct a simplicial complex as follows:*

- *Add a 0-simplex for each point in X .*
- *For $x_1, x_2 \in X$ add a 1-simplex between x_1, x_2 if $d(x_1, x_2) \leq r$.*
- *For $x_1, x_2, x_3 \in X$ add a 2-simplex with vertices x_1, x_2, x_3 if $d(x_1, x_2), d(x_1, x_3), d(x_2, x_3) \leq r$.*
- *...*
- *For $x_1, x_2, \dots, x_m \in X$, add a $(m-1)$ -simplex with vertices x_1, x_2, \dots, x_m if $d(x_i, x_j) \leq r$ for $0 \leq i, j \leq m$; that is, if all the points are within a distance of r from each other.*

The simplicial complex is called the Vietoris-Rips complex and is denoted $VR_r(X)$.

Definition 2. *We say that disjoint subsets $A_1 \dots A_k$ of vertices V in graph $G(V, E)$ are k connected components of the graph G if the following is true:*

1. *The vertices inside A_i are connected i.e. there exists a path between arbitrary two vertices $a, b \in A_i$, for every $i \in 1 \dots k$.*
2. *The sets of vertices A_1, \dots, A_k are disconnected i.e. there isn't an edges $e(a, b)$ between a pair of two points (a, b) such that $a \in A_i, b \in A_j, i \neq j$.*

Given the input of n m -dimensional points $X \in \mathbb{R}^m$ and a metric d (in our case the euclidean distance) we build a VR-cx for parameter $r = r_2$, where r_2 is the biggest r such that the VR-cx has three connected components. We find the parameter r_2 using a binary search on a set of all possible values of r . We then perform the classification, simply by saying that all the objects in some connected component belong to the same class. Seeing as we have three connected components in our VR-cx we will get three distinct classes.

It does not take much thought to see, that our classification model will produce the exact same results if we only use simplices of dimension 1, as it would, if we use all the simplices up to dimension h , where h is an arbitrary number in the range $h \in 1...m$. This follows from Definition 1.

Relation to single linkage clustering algorithm

Our intuition tells us, that the model we build using the VR-cx to classify the images, produces the same results as the well known single linkage clustering algorithm. In this section we aim to prove or at least give a strong intuition that this is indeed the case.

Both algorithms take a set X of n samples with m features as input. We can think of a sample x in the set X as a point in the m -dimensional space $x \in \mathbb{R}^m$. The algorithm then constructs a graph with points X as the vertices (V) and edges E . The k connected components in the constructed graph correspond to the classes of samples.

Single linkage algorithm. The algorithm starts with n connected components (no edges in the graph). In each step the algorithm chooses the two connected components that are closest to each according to some distance metric d (in our case the euclidean distance) and joins them into one by adding an edge between their closest two vertices.

Definition 3. The distance D between two connected components A and B is defined as the distance of the pair of vertices (one from A and one from B) that are closest to each other. More formally

$$D(A, B) = \min_{a \in A, b \in B} d(a, b).$$

The algorithm stops when there are only k connected components left.

Vietoris-Rips classification algorithm. The algorithm builds a (1-dimensional) Vietoris-Rips complex $V_r(X)$ with parameter r . We choose the biggest r such that the Vietoris-Rips complex $V_r(X)$ has k connected components.

To prove that the two algorithms indeed produce the same connected components we will first prove the next claim.

Claim 1. Let $G_{sl}(V, E_{sl})$ be a graph produced by the single linkage algorithm for finding k clusters and let d_{max} denote the distance between vertices in graph G_{sl} that were connected in the last iteration of the algorithm. Graph G_{sl} has connected components $A_1...A_k$. The graph $G_{vr}(V, E_{vr})$ induced by the Vietoris-Rips complex $VR_{d_{max}}(V)$ has the same connected components $A_1...A_k$.

Proof. We prove the Claim 1 by induction on the steps in the single linkage algorithm. We start with a set of vertices V . Let j denote the step (iteration) of the algorithm, e_j the edge added in j -th step and d_j its length. We claim that at each j the graph G_{sl}^j constructed by the algorithm up to that point, has the exact same connected components as G_{vr}^j , that is the graph induced by the Vietoris-Rips complex $VR_{d_j}(V)$.

Base case. For $j = 0$ this is obvious, since this is the initial state of the algorithm. Both graphs G_{sl}^0 and G_{vr}^0 consist only of vertices V . For $j = 1$ the algorithm adds the smallest edge e_1 out of all possible candidates and builds a graph G_{sl}^1 . Edge e_1 has length d_1 . It is obvious that $VR_{d_1}(V)$ will induce a graph G_{vr}^1 that will also only contain edge e_1 , since no other pairwise distance between vertices V is smaller.

Induction step. Here we show that if for some j our claim holds, it will also hold after another iteration of the algorithm i.e. for $j + 1$. In $(j + 1)$ -th iteration, the algorithm finds the edge e_{j+1} with length d_{j+1} and adds it to the graph. By the definition of the algorithm e_{j+1} is the smallest such edge that connects (joins) two separate connected components. This means that every other edge e' with length $d' < d_{j+1}$ would not join connected components, but would instead just connect two vertices, that are both already in the same connected component. From the definition of the Vietoris-Rips complex we can see that in the graph G_{vr}^{j+1} there will only be one new edge that will join two separate connected components, and that will be exactly edge e_{j+1} . All the other extra edges that will be added in G_{vr}^{j+1} , but do not appear in G_{sl}^{j+1} have length less than d_{j+1} and will therefore only connect vertices inside of already existing connected components of the graph G_{vr}^j . Since by our induction hypothesis graphs G_{sl}^j and G_{vr}^j had the same connected components and we joined two of the same connected components in both graphs, this means that the graphs G_{sl}^{j+1} and G_{vr}^{j+1} also have the same connected components.

We have proven that the graph G_{sl}^j constructed in j -th iteration of the single linkage algorithm indeed contains the same connected components as the graph G_{vr}^j induced by $VR_{d_j}(V)$ for an arbitrary j . This also proves Claim 1. \square

Using Claim 1 we see that the connected components in G_{vr} and G_{sl} are indeed the same. We need to take into account that the Vietoris-Rips algorithm takes the biggest such r , so that the graph has k connected components, so $r > d_{max}$. But we can quickly see that the extra edges in the graph induced by $V_r(V_{sl})$ will not change the connected components. After all we already have k connected components in G_{vr} . To join any two together would mean a violation of a fundamental rule of the algorithm.

Computational complexity

Let us now consider the computational complexity of our model. Since we are only interested in Vietoris-Rips complexes VR_r with simplices of dimension 1, the simplest approach to construct such complex requires us to check the distance between every pair of vertices $(x_1, x_2) \in X \times X$, adding such pair to the final complex if the distance $d(x_1, x_2) \leq r$. In worst case the algorithm would have to return all distinct pairs of vertices, meaning construction of $VR_r(X)$ requires $O(n^2)$ time and consumes $O(n^2)$ space, where n is the number of vertices in X .

The problem is that we don't know the appropriate value for the parameter r . Recall that we are interested in finding biggest r , such that the Vietoris-Rips complex VR_r has as many connected components as there are distinct classes of images. Let r_{max} denote the largest distance between two vertices from X . Note that r we are looking for will always be bounded on the interval $[0, r_{max}]$, and it will furthermore be exactly one of the distances between some pair of vertices. Thus we only have n^2 different possible values of r to check, and if we sort them by size and use binary search to find the right one, we can do it in $O(n^2 \log n^2) = O(n^2 \log n)$ time and $O(n^2)$ space. To count the number of connected components obtained with every different VR complexes, we can use a union-find algorithm, which at worst adds $O(n^2)$ time to each run.

With this the final time complexity of our approach is $O(n^2 \log n^2 + (n^2 + n^2) \log n^2) = O(n^2 \log n)$ using $O(n^2)$ space. Contrast this with the computational complexity of single-linkage clustering, which with a clever implementation can in optimal case produce solution in $O(n^2)$ time and $O(n)$ space [1].

Results

Cup, paper, pen image set

Effect of preprocessing We were interested in observing the effect of the preprocessing methods on the samples of our dataset. For this purpose we visualized the data-points using multidimensional scaling

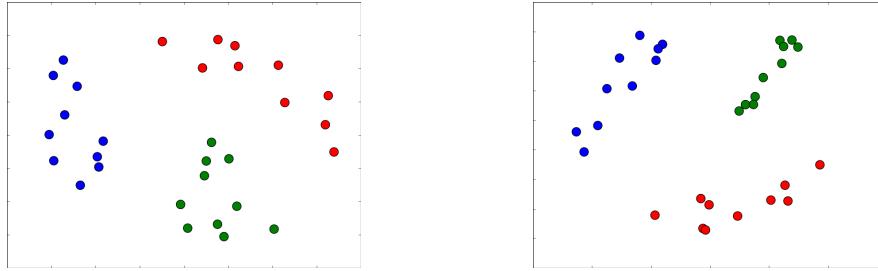


Figure 3: MDS plots visualizing cup (blue), paper (red), and pen (green) samples from the dataset. Left figure is the original dataset. Right figure shows the scaled and PCA transformed dataset. Euclidean distance metric was used.

technique (MDS) [?]. MDS is a means of visualizing the level of similarity of individual cases from N -dimensional space, to a lower, M dimensional space (in our case: $M = 2$), while preserving between-object distances as good as possible. Therefore MDS plots are well suited for giving intuition how separable data samples of different classes are. The obtained plots are shown in figure 3. We notice that in both cases the points can be clearly separated into three clusters. We can also observe that points within clusters seem to lie closer together in the preprocessed plot than in the non preprocessed one. It should be mentioned that the position of points in a MDS plot is randomized, and that only distances between points should be observed. Therefore the reader should not confuse the different positions of each of the clusters in the plots as an effect of the preprocessing, since the reason for it is merely a property of the MDS method.

Image classification. In this section we present the results of Vietoris-Rips classification algorithm described in section . TODO: Why is this reference not shown?

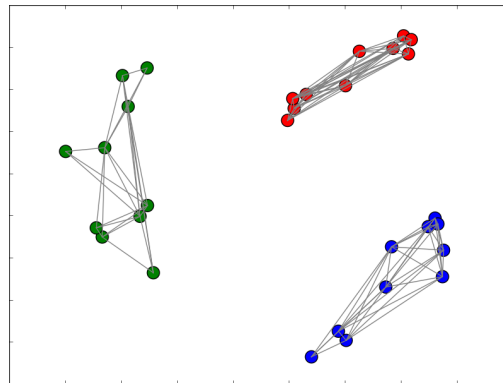


Figure 4: Visualisation of the VR-cx on the cup, paper, and pen dataset.

Initially we tested if the model is able to distinguish the already mentioned images of cup, paper and pen. The images were preprocessed, as explained before, with scaling and a PCA transformation. The VR-cx obtained is shown in fig 4. The model had no problems distinguishing between the three classes of images. Each one of the three connected components obtained, was belonging to one of the three classes of images. This is not surprising since looking already mentioned MDS plots 3, the three different classes of points are clearly separable.

Principal simplices are simplices which are not faces of any bigger complex. In the built VR-cx we observed 14 principal simplices, which are drawn in figure 5. We can notice that each of the principal

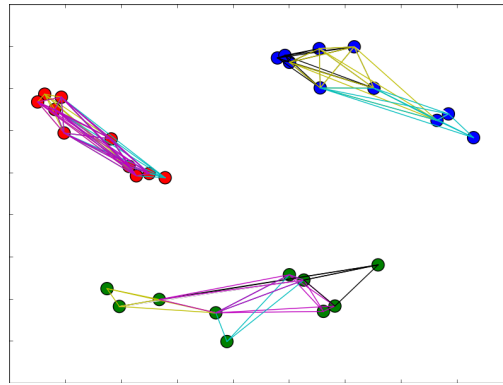


Figure 5: Visualisation of the principal simplicies obtained by VR-cx on the cup, paper, and pen dataset.

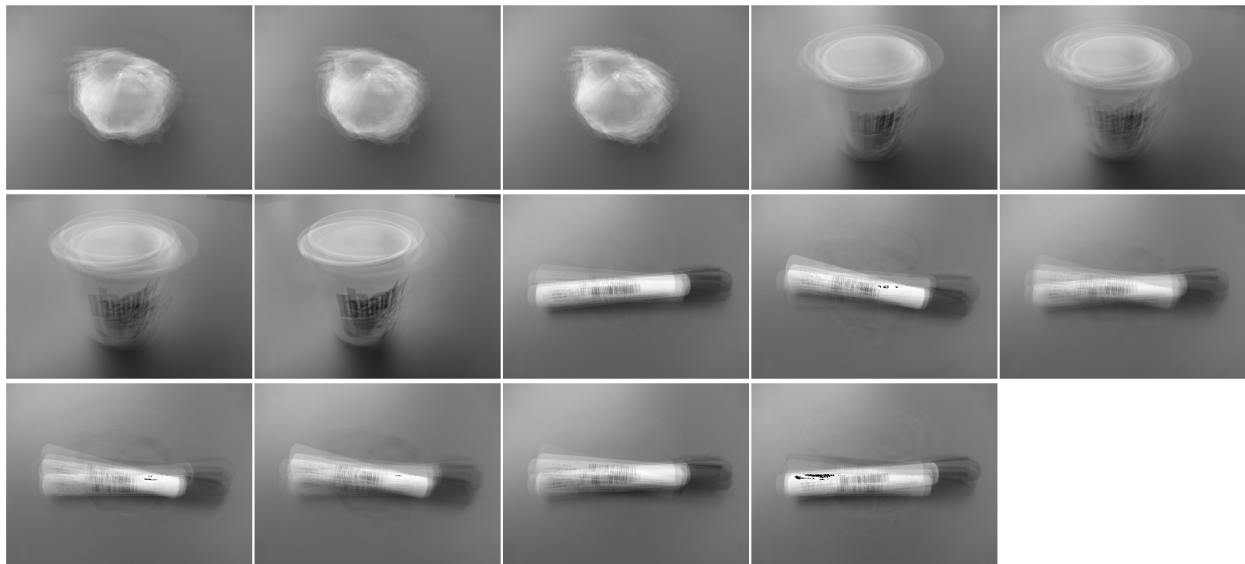


Figure 6: Images representing inverse-transformations of baycenter points of each of the principal simplicies ordered by decreasing sizes of simplicies.

simplicies is connecting points with same classes. In 6 the barycenter images of each of the simplicies is shown. **TODO: Describe results**

TODO: [A] Explain that simplicies of larger dimension are just about useless

Intercluster edges **TODO: [A]** Pick one (or several) edges in V_1 which is not in V_2 and reconstruct the image corresponding to the midpoint. This should somehow capture the difference between the two components connected by this edge. Does it?

General classification

TODO: [R] Explain we get same results as with S-L clustering, and with that the same problems as S-L clustering.

TODO: [R] Show bad results on our first dataset.

TODO: [A] Extra tests: iris, ...

Modifying algorithm to prefer simplices of larger dimension. Because we proved that our basic model that only considers simplices of dimension 1 produces identical results as single-linkage clustering algorithm, and because VR cx can produce simplices of higher dimensions, we wondered if we could somehow use those to increase robustness of our algorithm. The main idea is that if some vertices are present in a simplex of higher dimension, they are more likely to be a part of the same class (since all the points are bunched up closer together). Unfortunately this turned out to be useless in practice, as even simplices of higher dimension are still built based on some distance metric between points, and if two clusters of different classes are close enough together, they may form a simplex of very large dimension resulting in a incorrect classification. The idea also would not work too well on thin, long clusters as seen in figure 3. On top of everything, this would also add a second parameter (dimension in addition to r), that the model has to balance out to produce desired number of clusters, and we would have to construct simplices of dimensions ≥ 1 at each run of VR algorithm, further increasing the computational complexity with potentially no benefit.

Summary

TODO: [A] Brief summary

TODO: [A] Further work (if there even is any), what we didn't explore.

Workload

TODO: Write here about your work on the project and I ([A]) will then summarize it into a paragraph or two.

- P: Prepared Preprocessing class (StandardScaler, PCA), prepared VR-cx, added MDS visualizations, tested different scalers.
- R:
- A: Implemented methods for loading/saving images, preparing the dataset for use in the model, implemented the model, converted original images to grayscale and resized them,

A more detailed overview of each member's workload can be seen at github repository of our project: https://github.com/ptrus/vietoris_rips_image_classification/commits/master.

References

- [1] John C Gower and GJS Ross. Minimum spanning trees and single linkage cluster analysis. *Applied statistics*, pages 54–64, 1969.