

Introduction to MAPPER

Leyda Almodóvar

You will find Mapper and instructions to download it and install it here:
<http://danifold.net/mapper>

Or see page 2 of this document

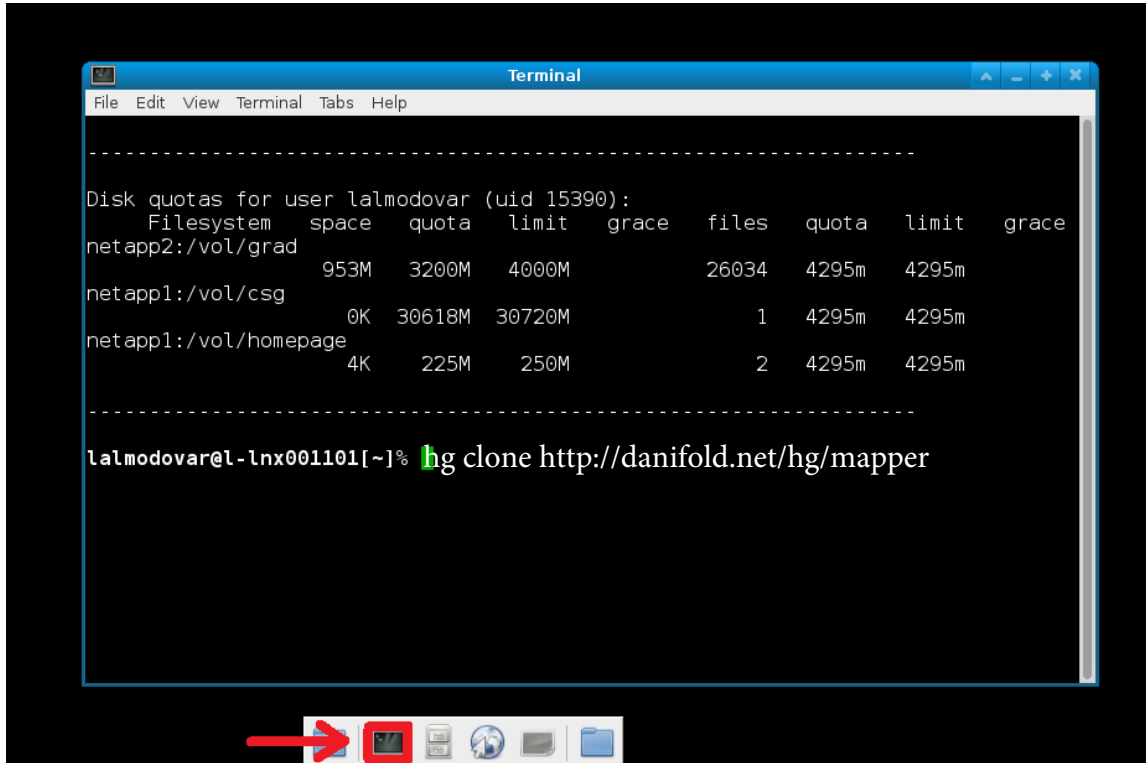


Make sure to look at <http://danifold.net/mapper/installation/index.html> for installing instructions and troubleshooting tips.

It is important that your data contains only numbers. Remove dollar signs, commas, NA, NaN and headers.

In order to download Mapper open up a terminal (found on the bottom of the screen):

Figure 1: Open up a terminal

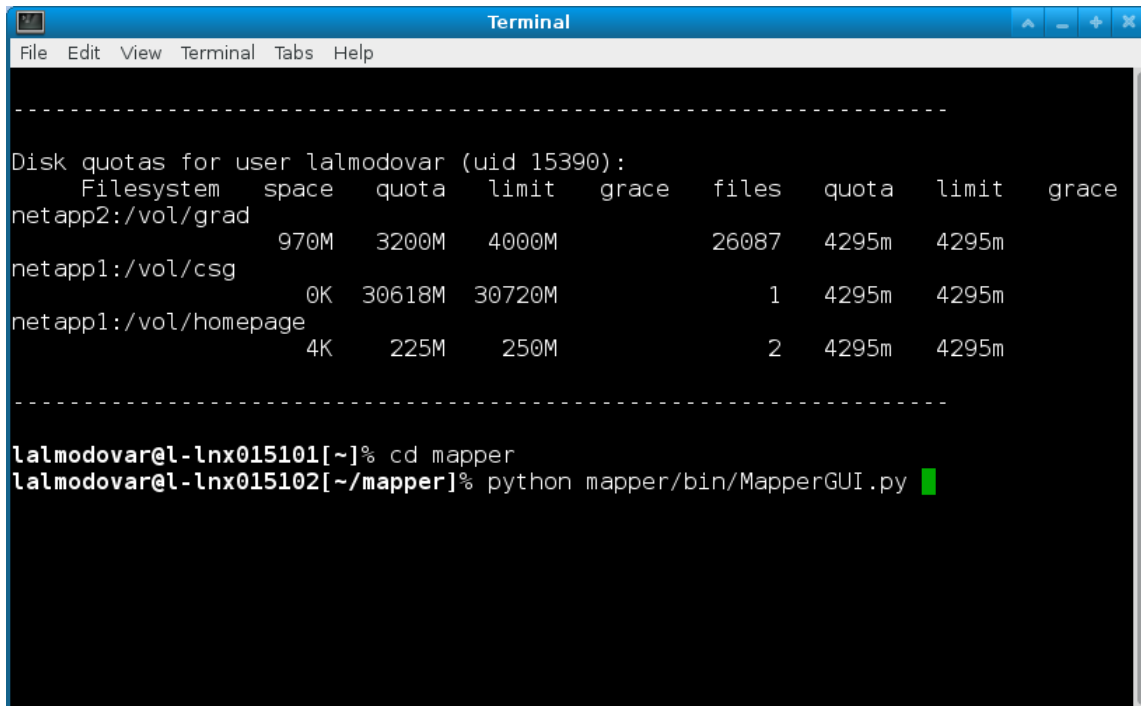


Download: Type the following in a terminal window:
hg clone http://danifold.net/hg/mapper

In order to open Mapper you could double click the folder named 'mapper' and then double click the folder named 'bin' and then double click the file named 'MapperGUI.py'.

Alternatively, you could type **cd mapper** then press enter, then type **python mapper/bin/MapperGUI.py** and press enter.

Figure 2: In order to open the GUI from the terminal:



```
-----  
Disk quotas for user lalmodovar (uid 15390):  
  Filesystem  space   quota   limit   grace   files   quota   limit   grace  
netapp2:/vol/grad      970M   3200M   4000M           26087   4295m   4295m  
netapp1:/vol/csg              0K  30618M  30720M             1   4295m   4295m  
netapp1:/vol/homepage    4K     225M    250M             2   4295m   4295m  
-----  
lalmodovar@l-nx015101[~]% cd mapper  
lalmodovar@l-nx015102[~/mapper]% python mapper/bin/MapperGUI.py █
```

Figure 3: Once you press enter the GUI should appear and it looks like this. The next step is to load the data.

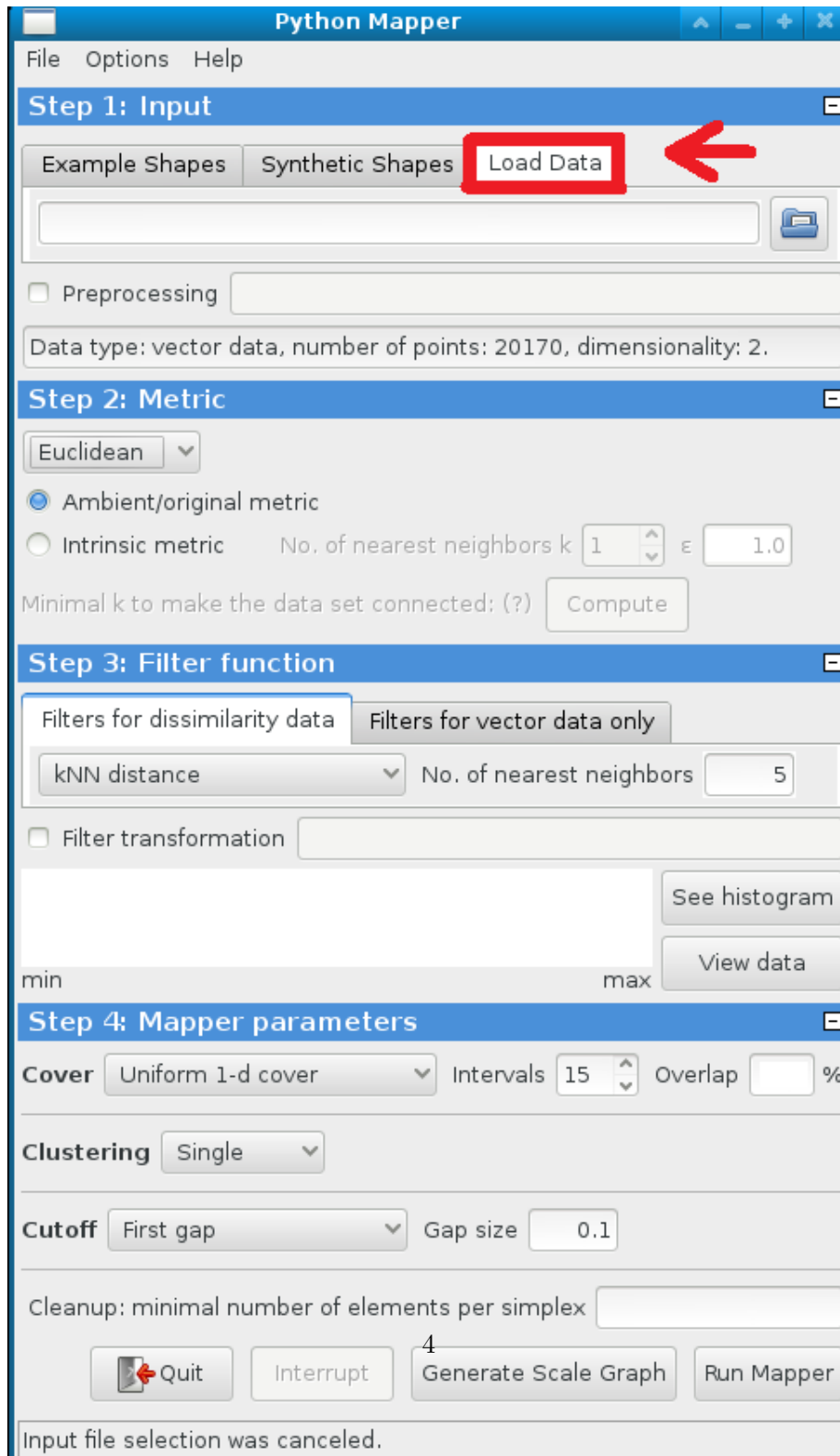


Figure 4: Browse for data file

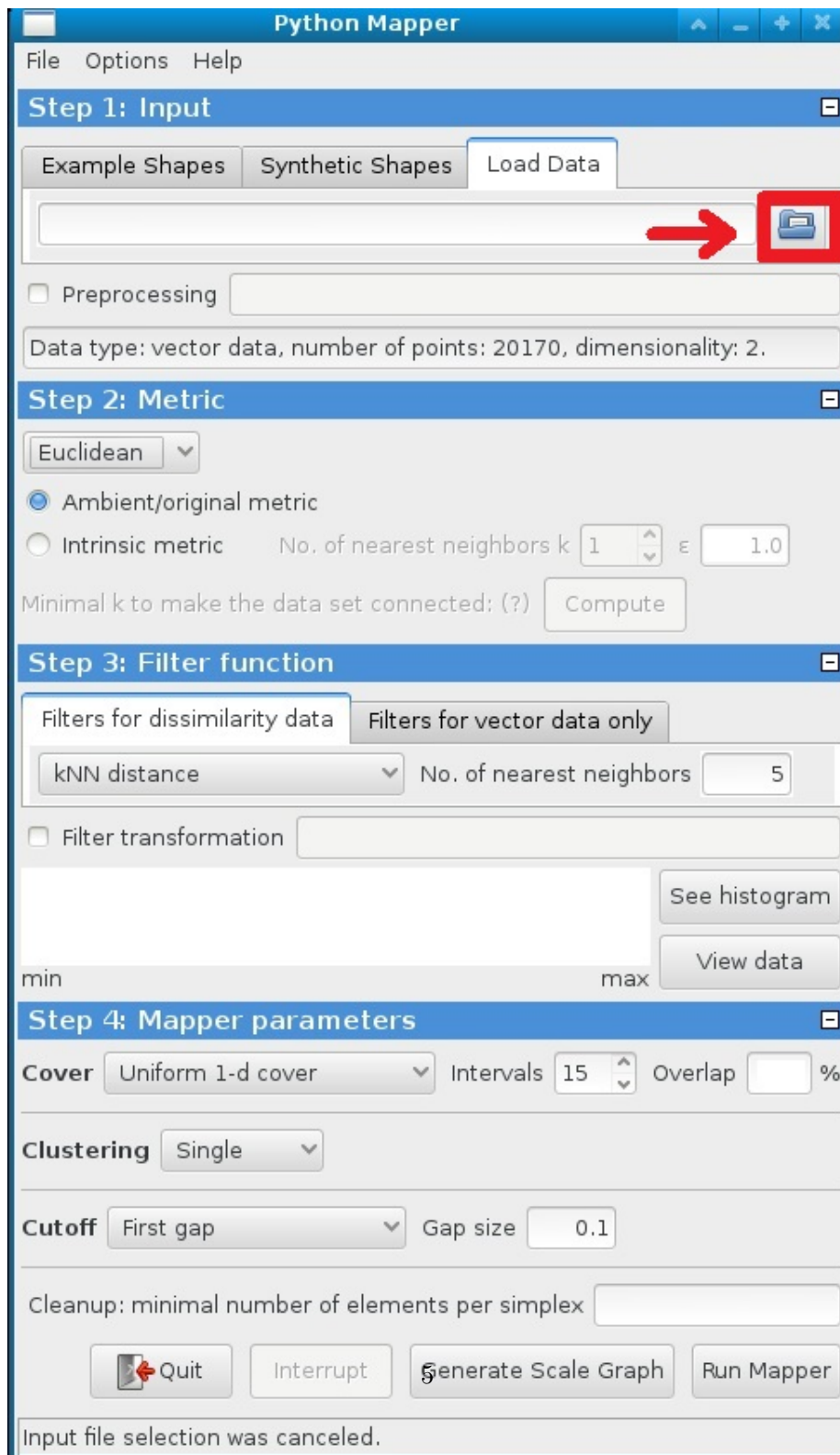


Figure 5: Mapper tells you the number of data points and dimensionality of your data set once you select the file

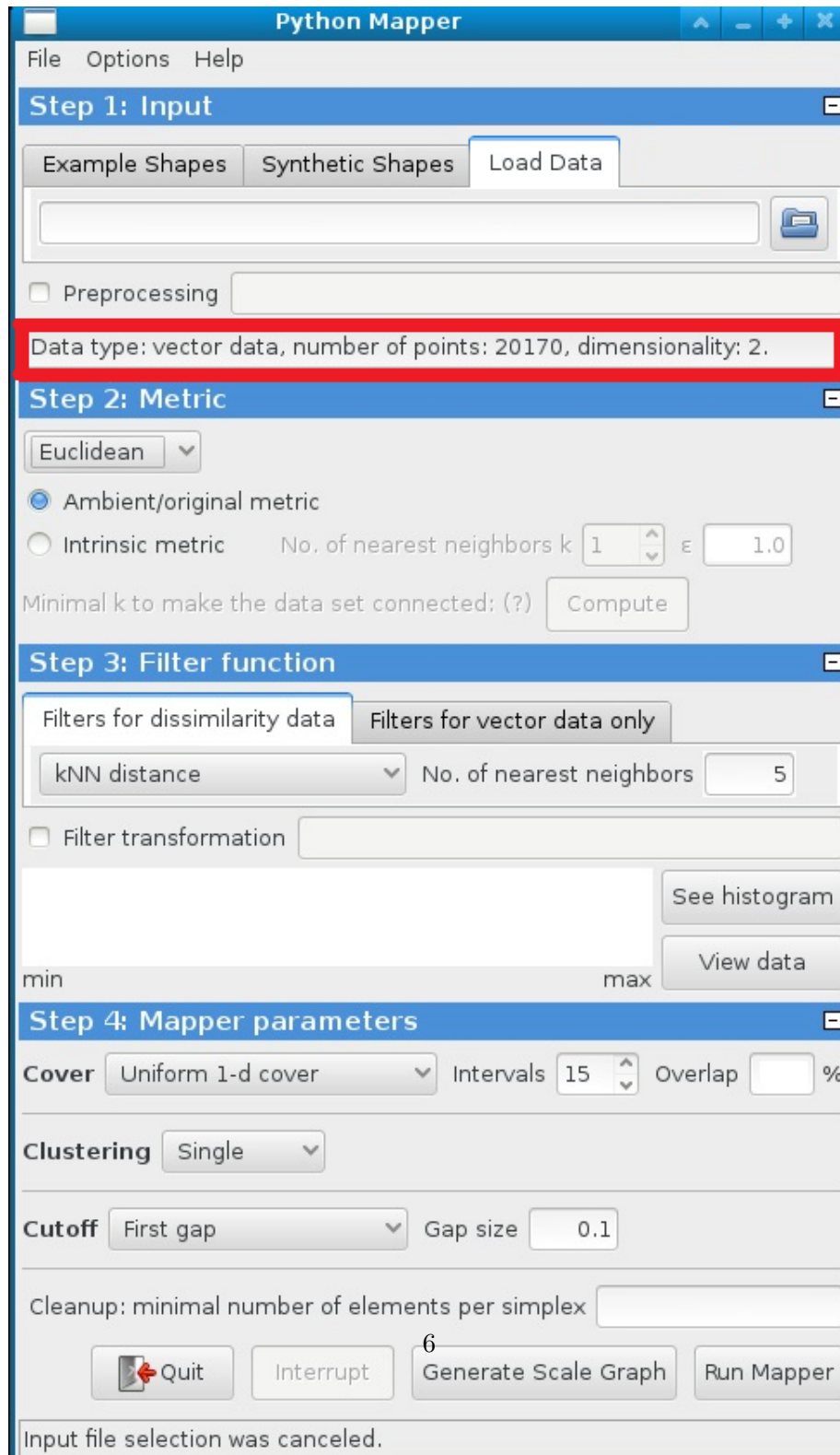


Figure 6: Choose a metric: Euclidean, Minkowski or Chebychev

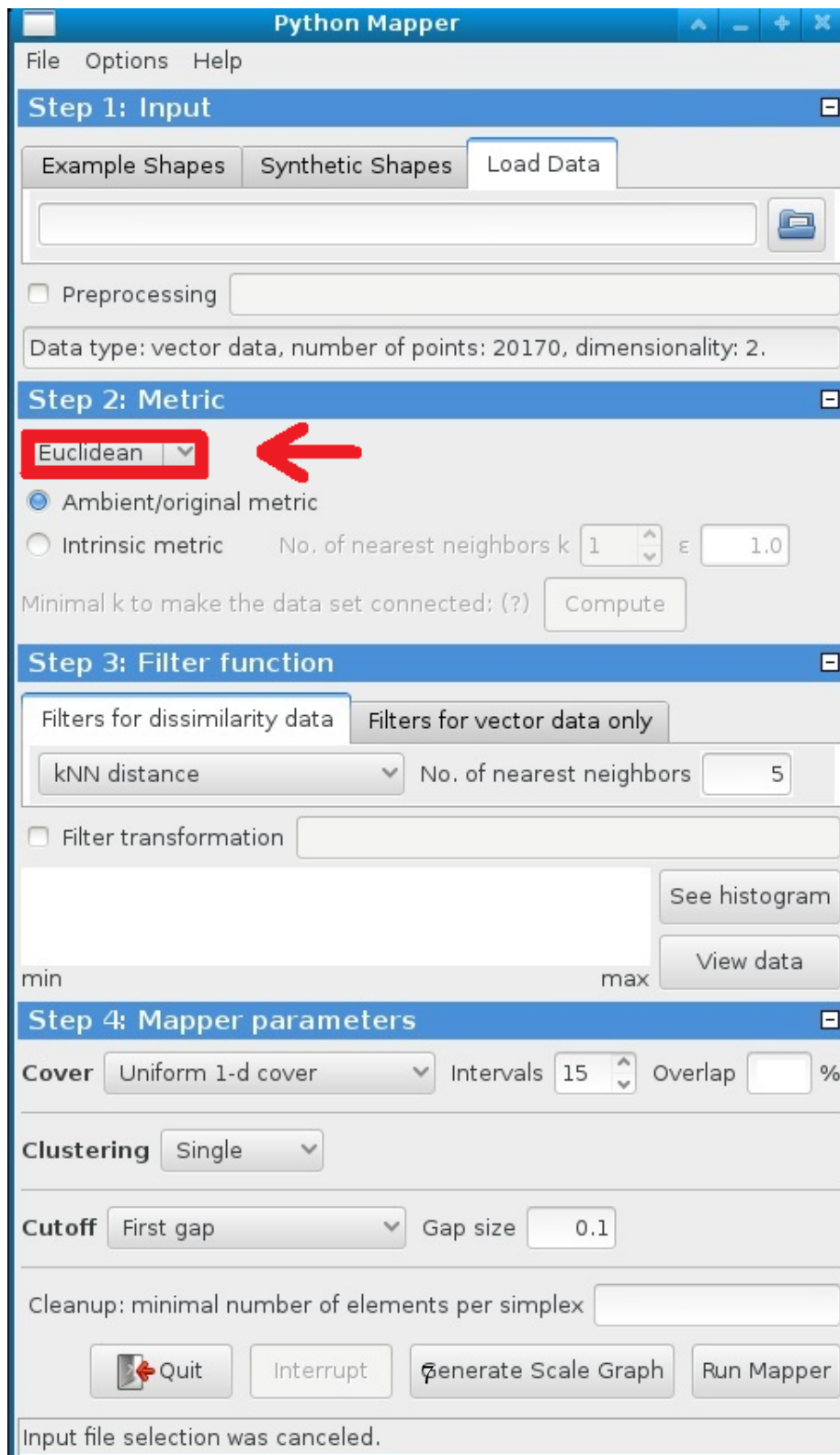


Figure 7: Choose a filter function: Eccentricity, kNN distance, Distance to a measure, Density (Gaussian Kernel), Graph Laplacian or Distance matrix eigenvector

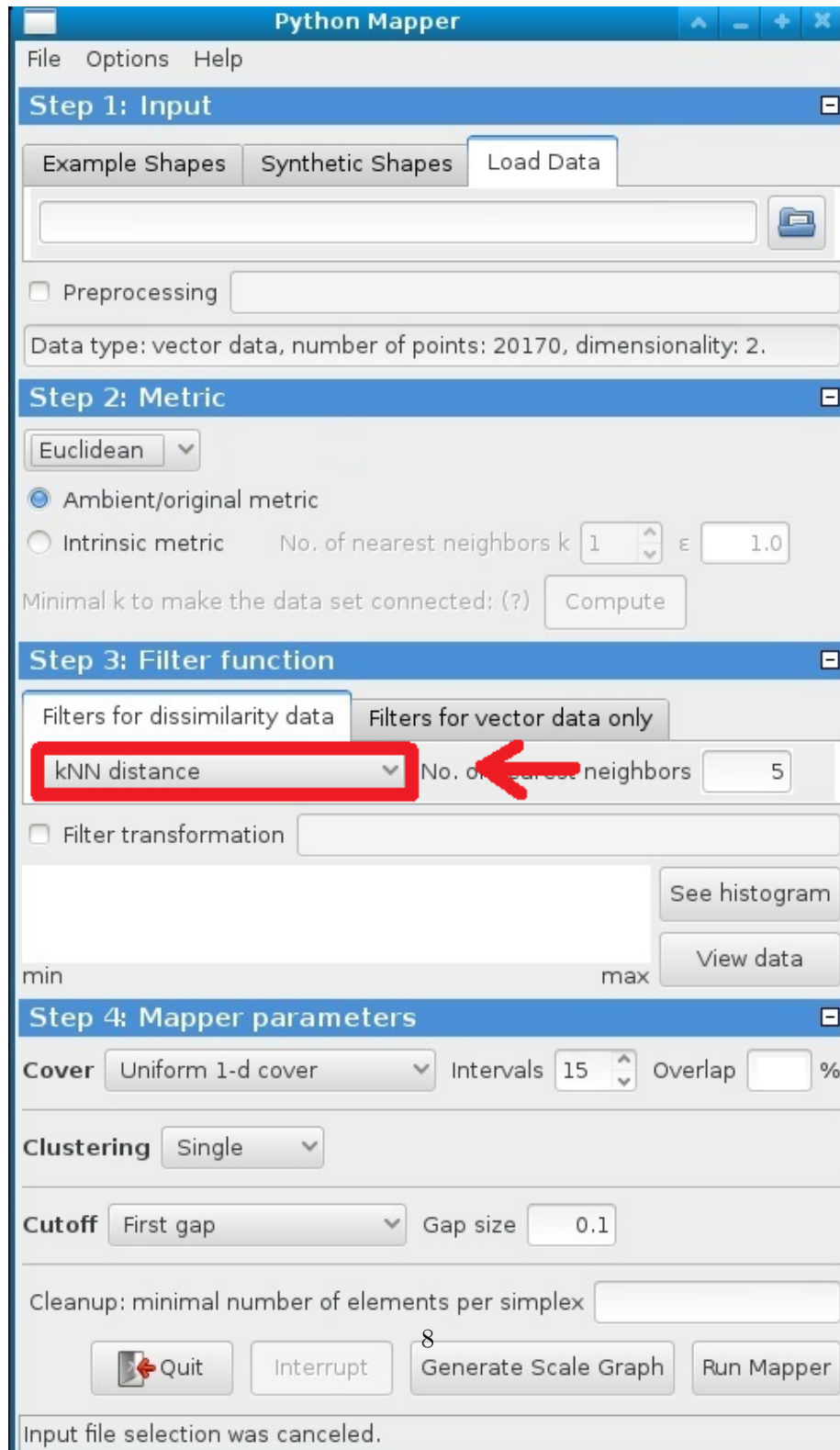


Figure 8: View data (as long as it is 1-dimensional, 2-dimensional or 3-dimensional)

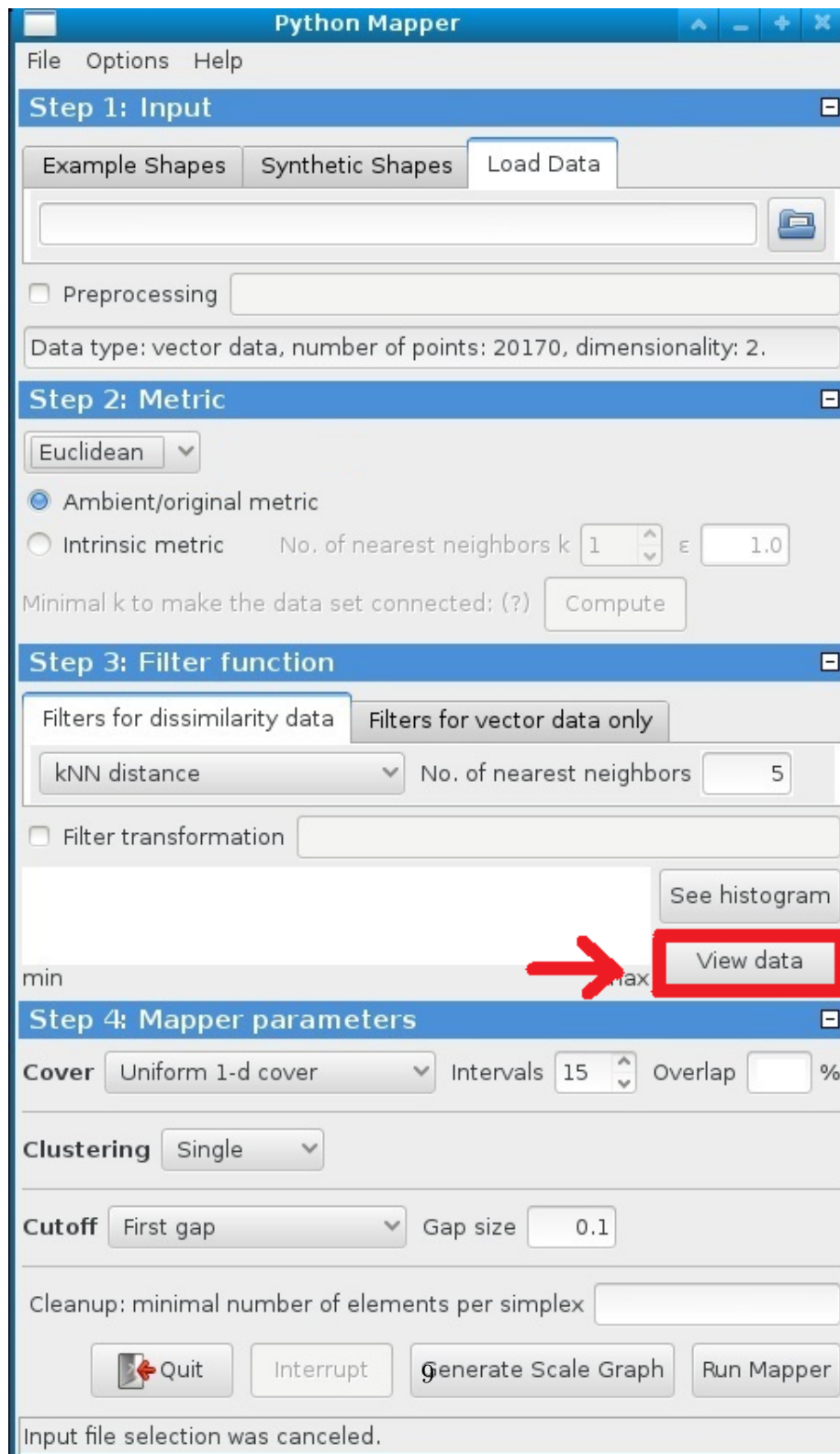


Figure 9: Choose the type of cover, number of intervals and percentage of overlap between successive intervals

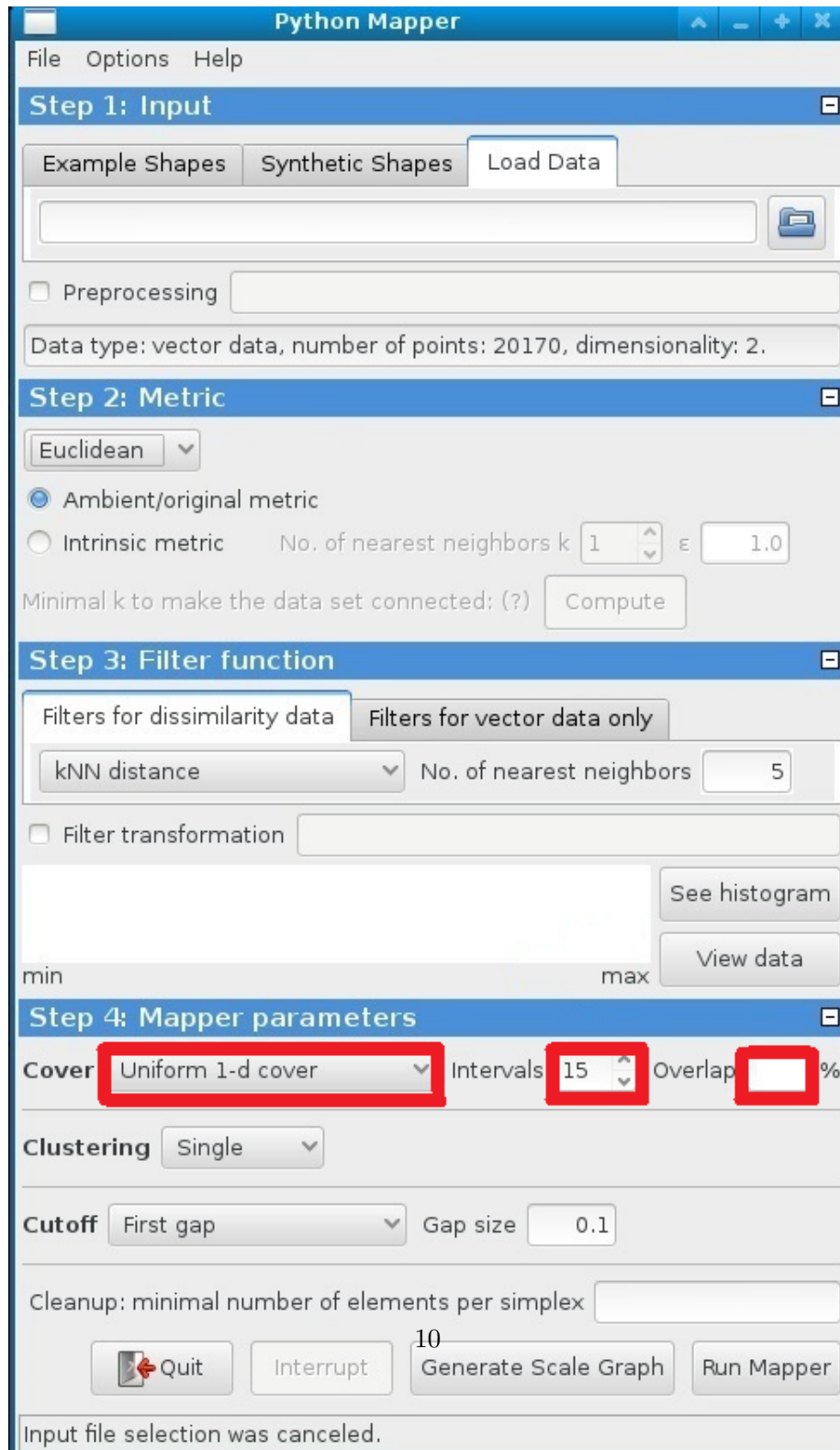


Figure 10: Choose clustering algorithm: Single, Complete, Average, Weighted, Median, Centroid, Ward

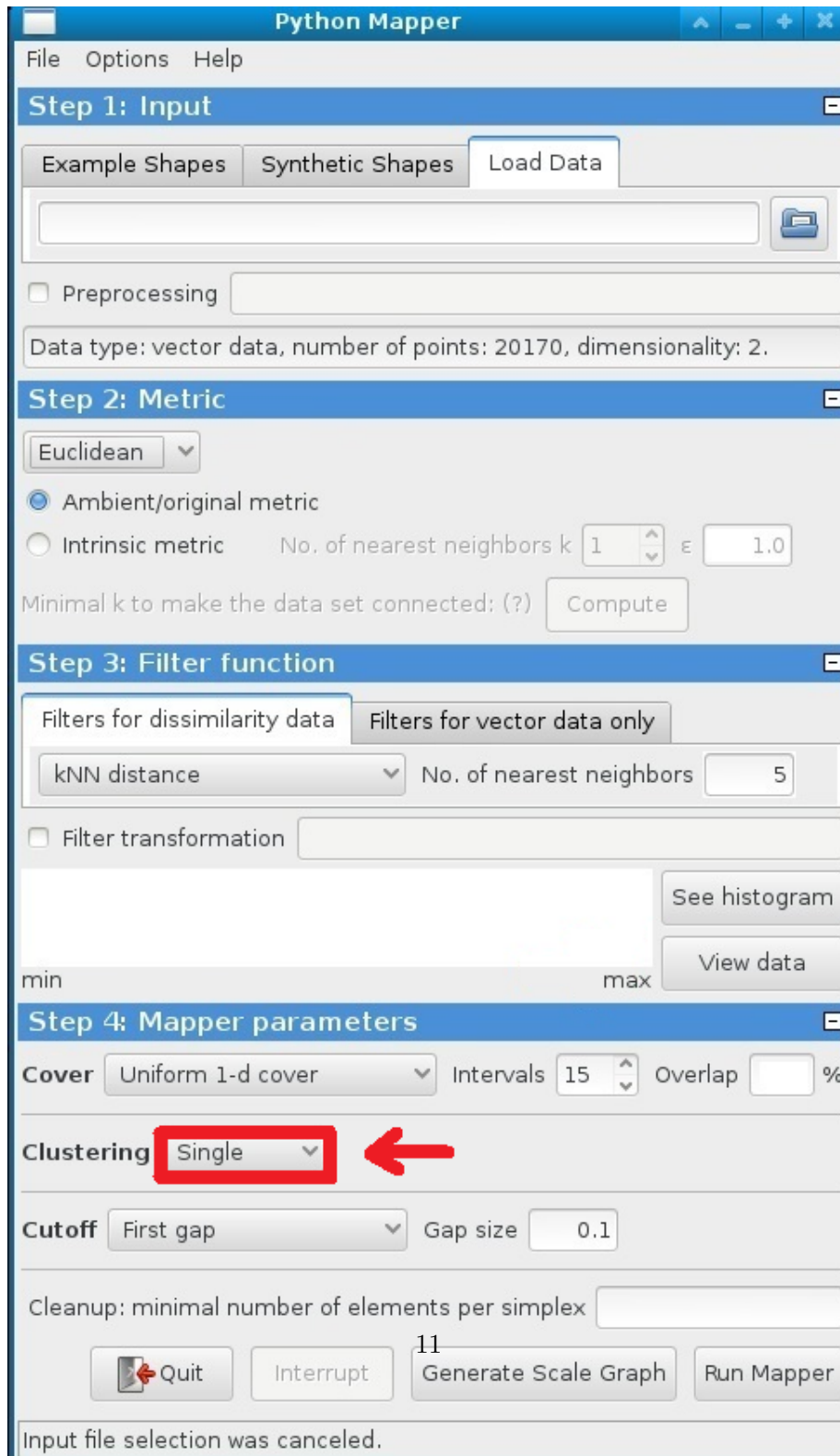


Figure 11: Run MAPPER

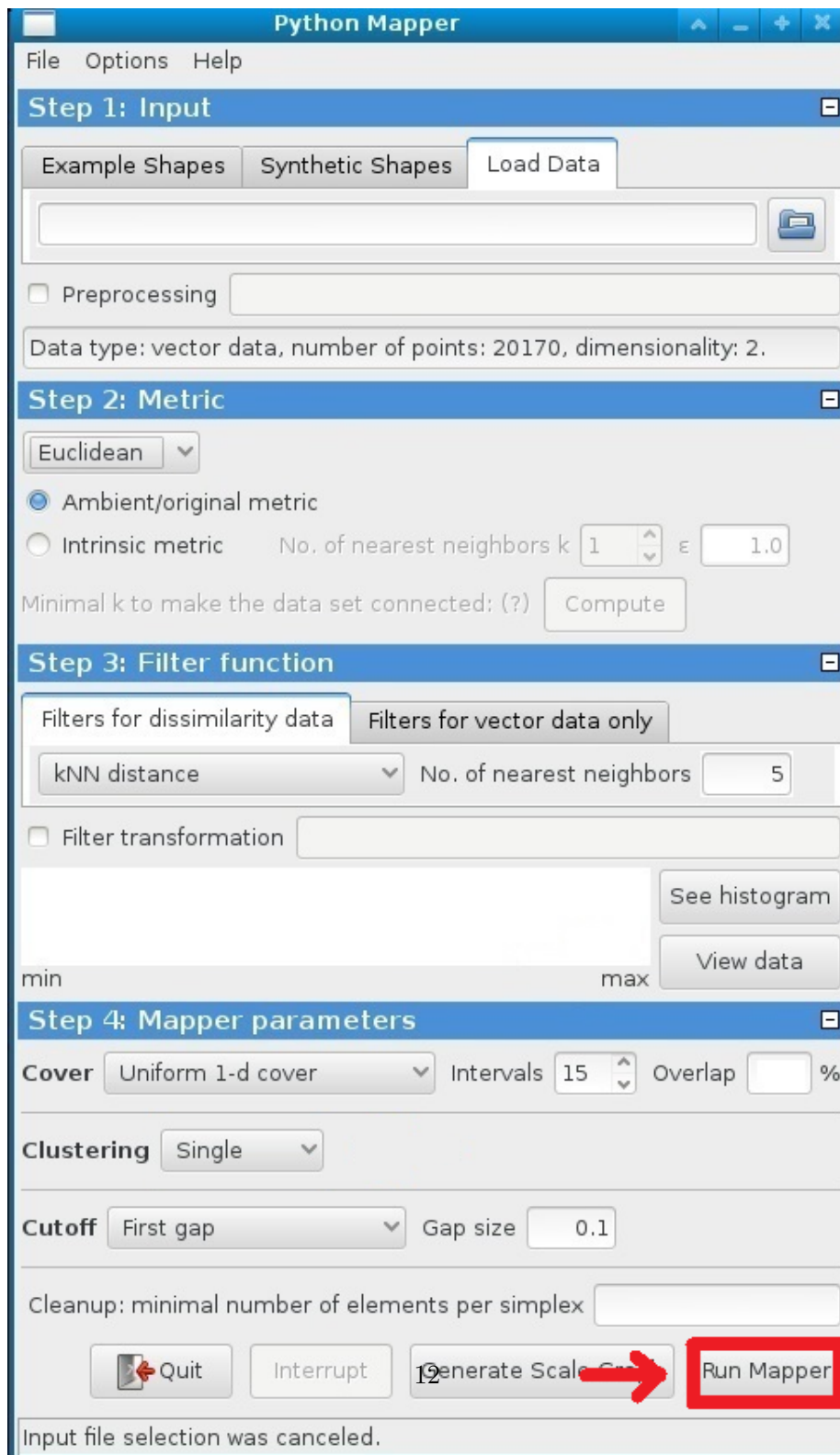
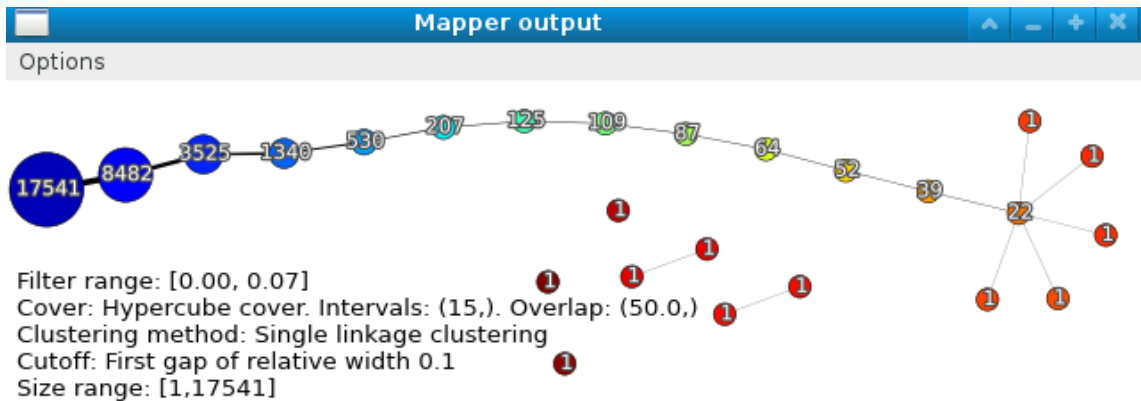


Figure 12: Sample Output



In order to run Mapper you must make certain choices.

1. Metric

- Euclidean*: The *Euclidean* metric between two points $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and $\mathbf{y} = (y_1, y_2, \dots, y_n)$ is given by $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{\sum_{i=1}^n |x_i - y_i|^2}$.
- Minkowski*: The *Minkowski* metric of order p between two points $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and $\mathbf{y} = (y_1, y_2, \dots, y_n)$ is defined as $d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n (|x_i - y_i|^p)^{1/p}$.
- Chebyshev*: The *Chebyshev* metric between two points $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and $\mathbf{y} = (y_1, y_2, \dots, y_n)$ is defined as $d(\mathbf{x}, \mathbf{y}) = \max_i |x_i - y_i|$.

2. Filter function

A filter function is a function on the data set, $f : X \rightarrow \mathbb{R}^k$. The Mapper algorithm supports general, vector-valued functions, while the GUI is restricted to real-valued functions (the case $k = 1$) for simplicity.

- Eccentricity
- kNN distance
- Distance to a measure
- Density, Gaussian Kernel
- Graph Laplacian
- Distance matrix eigenvector

3. Type of cover

- Uniform 1-d cover

- (b) Balanced 1-d cover
- (c) Subrange decomposition

4. Clustering algorithm

- (a) *Single*: The distance between two clusters is defined as that of the closest pair of individuals, where it only considers pairs consisting of elements from different clusters: $d(r, s) = \min(\text{dist}(x_{ri}, x_{sj}), i \in (1, \dots, n_r), j \in (1, \dots, n_s))$.
- (b) *Complete*: The distance between two clusters is defined as that of the most distant pair of individuals, where it only considers pairs consisting of individuals from different clusters: $d(r, s) = \max(\text{dist}(x_{ri}, x_{sj}), i \in (1, \dots, n_r), j \in (1, \dots, n_s))$.
- (c) *Average*: The distance between two clusters is defined as the average of the distance between all pairs of individuals that are made up of one individual from each cluster: $d(r, s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} \text{dist}(x_{ri}, x_{sj})$.
- (d) *Weighted*: The distance between two clusters is defined as the weighted average of the distance between all pairs of individuals that are made up of one individual from each cluster. It uses a recursive definition for the distance between two clusters: $d(r, s) = \frac{(d(p,s)+d(q,s))}{2}$.
- (e) *Median*: The distance between two clusters is defined as $d(r, s) = \|\tilde{x}_j - \tilde{x}_s\|_2$ where \tilde{x}_r and \tilde{x}_s are weighted centroids for the clusters r and s and \tilde{x}_r is defined recursively as $\tilde{x}_r = \frac{1}{2}(\tilde{x}_p + \tilde{x}_q)$.
- (f) *Centroid*: The distance used is the Squared Euclidean distance between centroids $d(r, s) = \|\tilde{x}_j - \tilde{x}_s\|_2$ where $\tilde{x}_r = \frac{1}{n_r} \sum_{i=1}^{n_r} x_{ri}$.
- (g) *Ward*: The distance used is $d(r, s) = \sqrt{\frac{2n_r n_s}{(n_r + n_s)}} \|\tilde{x}_r - \tilde{x}_s\|_2$. The distance is defined as the incremental sum of squares, that is, the increase in the total within-cluster sum of squares as a result of joining two clusters. The within-cluster sum of squares is defined as the sum of the squares of the distances between all objects in the cluster and the centroid of the cluster.

To find more information about the covers provided by Mapper:
<http://danifold.net/mapper/cover.html>

To find more information about the filters provided by Mapper:
<http://danifold.net/mapper/filters.html>