

Neural and Evolutionary Computation

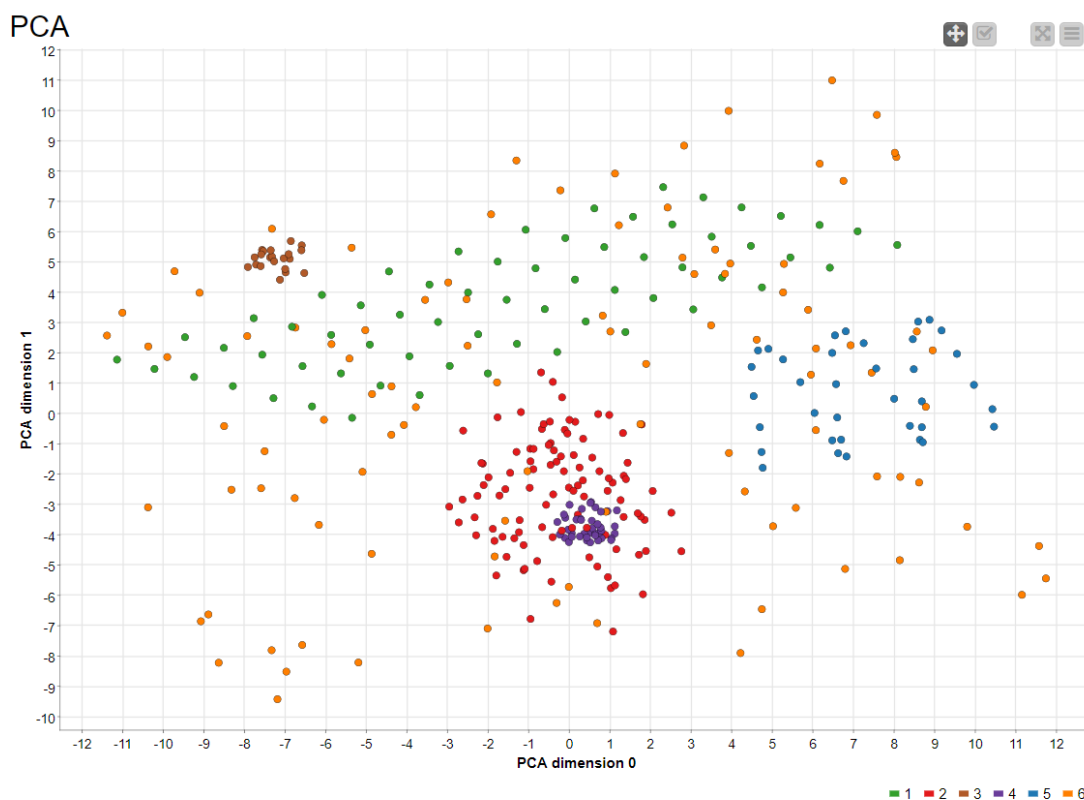
A3: Unsupervised learning with PCA, t-SNE, k-means and SOM

Introduction:

The goal of this assignment is to compare four unsupervised learning methods: PCA, t-SNE, K-means, and SOM. All methods will be compared with the same dataset. The dataset contains 360 patterns of four variables and one class column. In this report, we will discuss the different results after giving a short explanation of the implementation. As done in A2, I will use Knime for the implementation, despite for SOM which will be done with R. Please find attached to this report the results and the scripts.

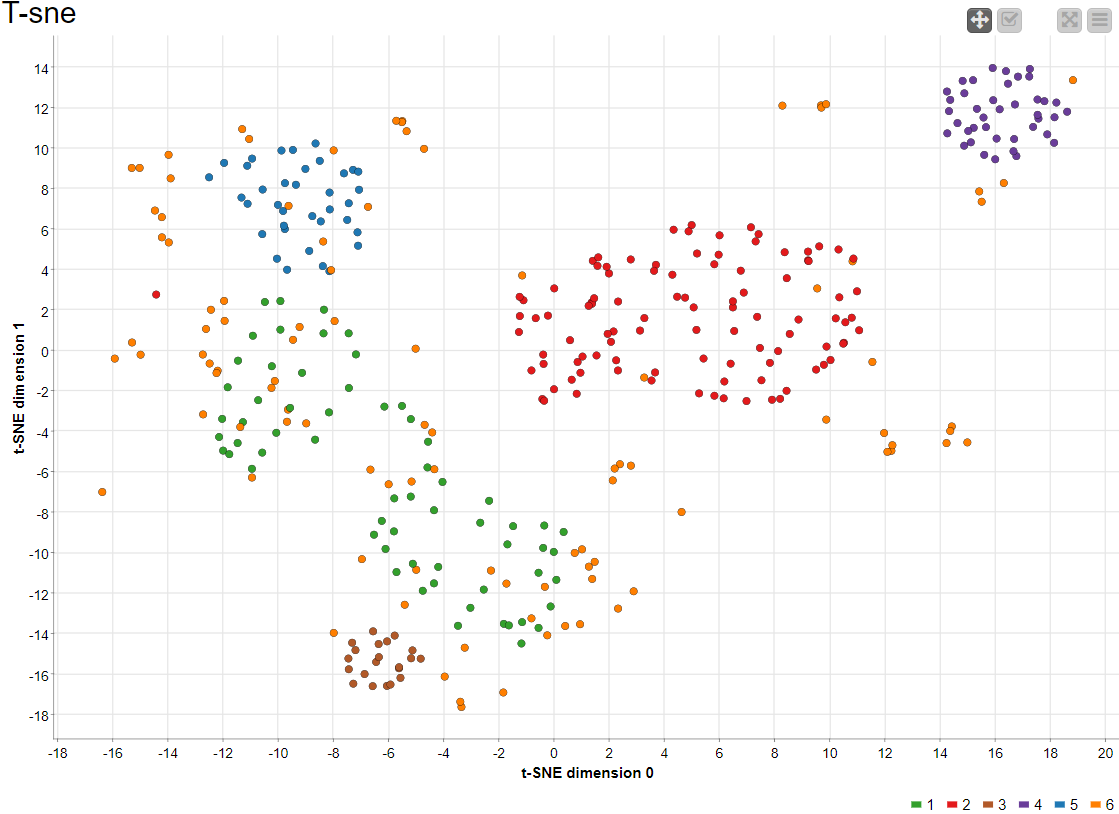
Description of the implementation:

PCA : Implementing PCA is quite fast. Knime nodes for PCA do not offer any configuration. At the end, we get this result.



t-SNE : Same as PCA we plot the dataset in two dimensions.

T-sne

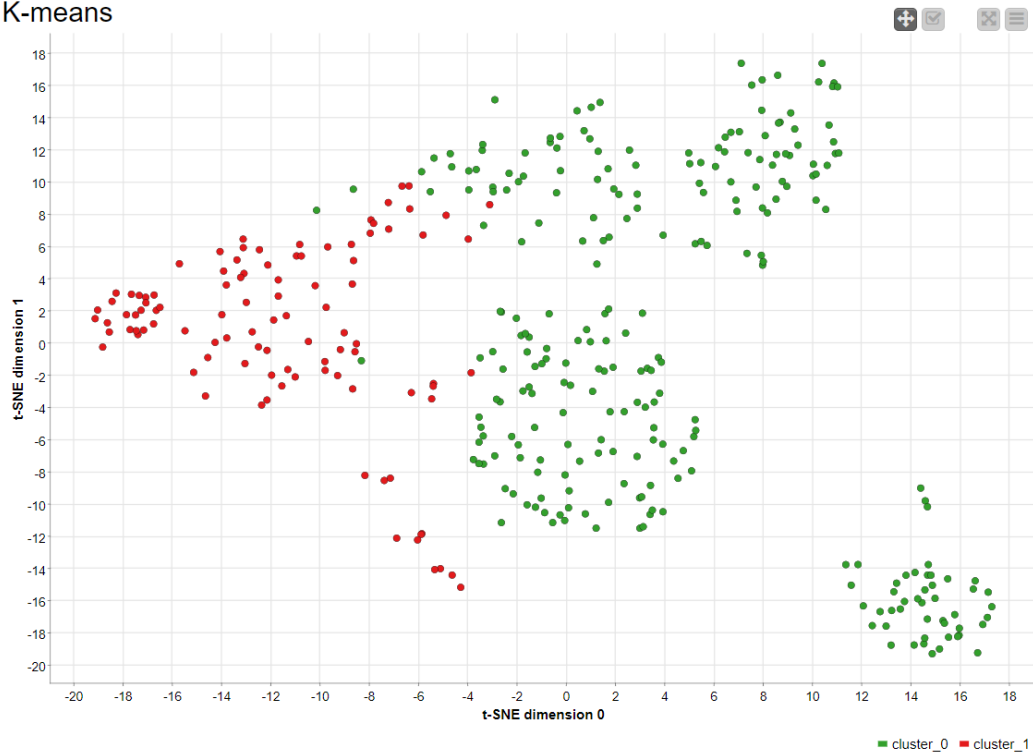


Using these parameters :

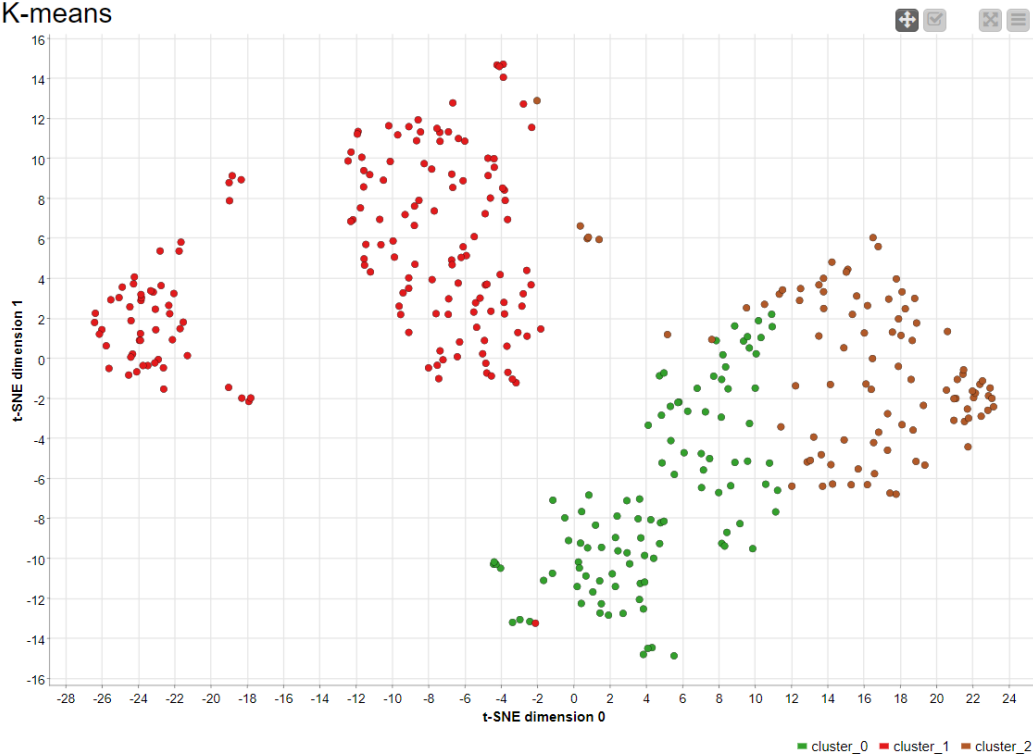
Dimension(s) to reduce to	<input type="text" value="2"/>
Iterations	<input type="text" value="1 000"/>
Theta	<input type="text" value="0,5"/>
Perplexity	<input type="text" value="30,0"/>
Number of threads	<input type="text" value="8"/>

K-means : For k-means we did several plots by varying the number of clusters from 2 to 7. The plots were made using t-SNE to reduce dimensions to 2 .We will discuss the results in next part.

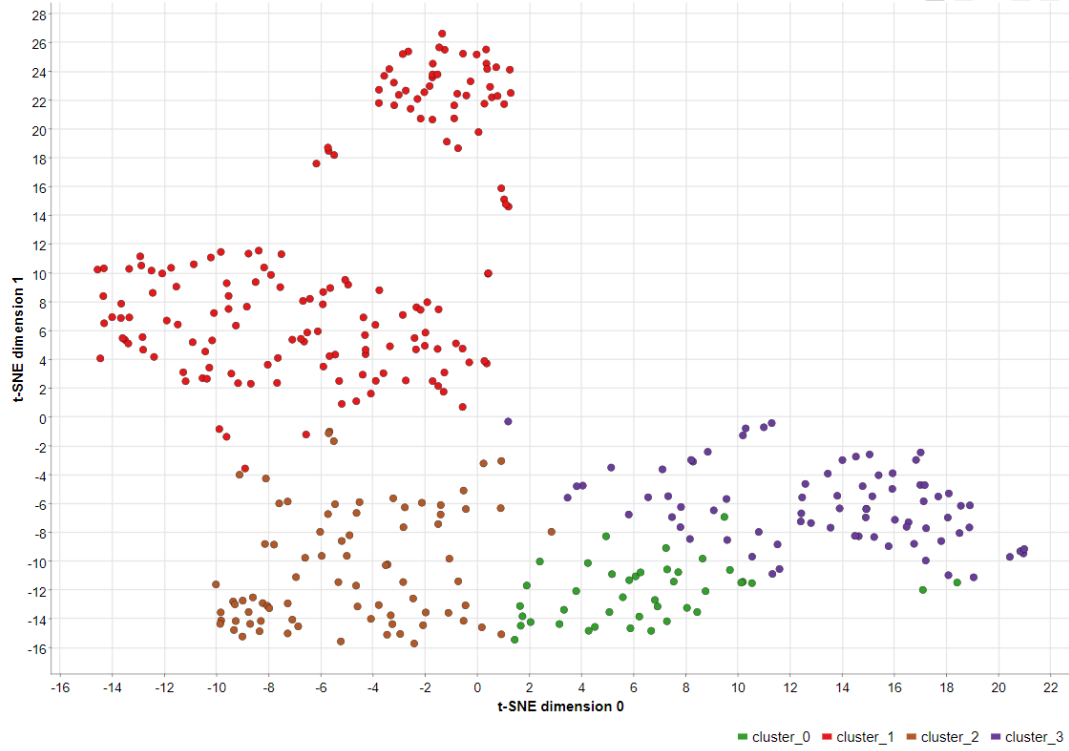
K-means



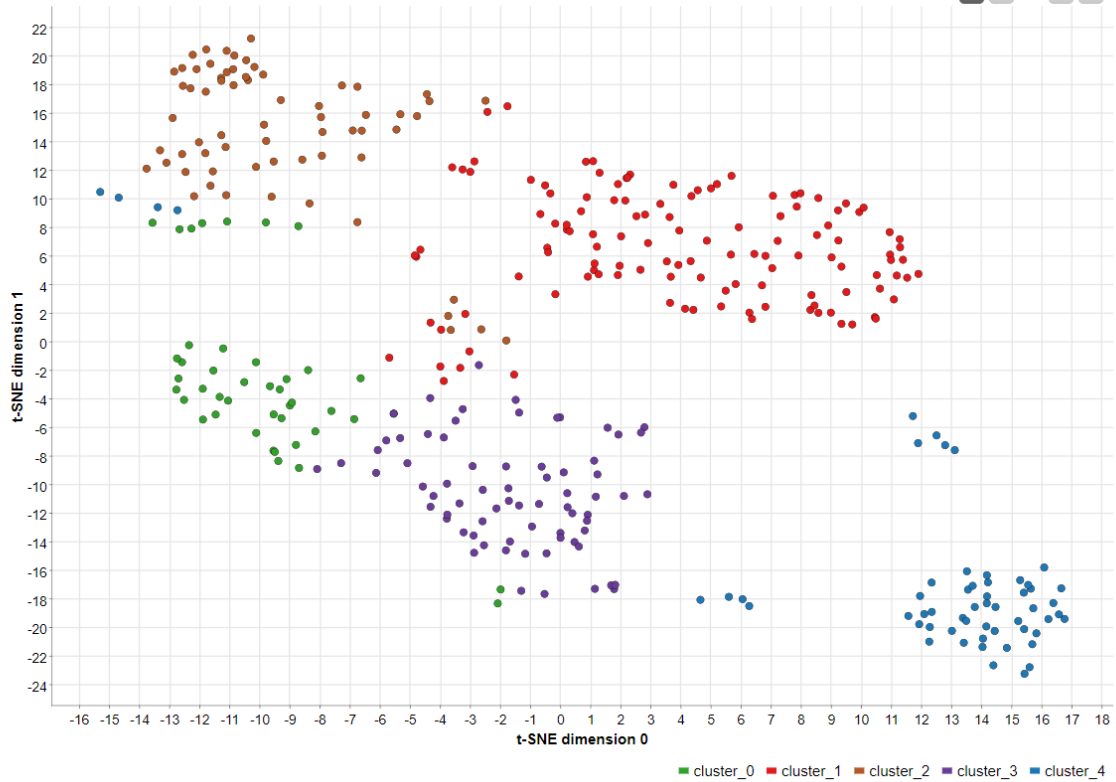
K-means



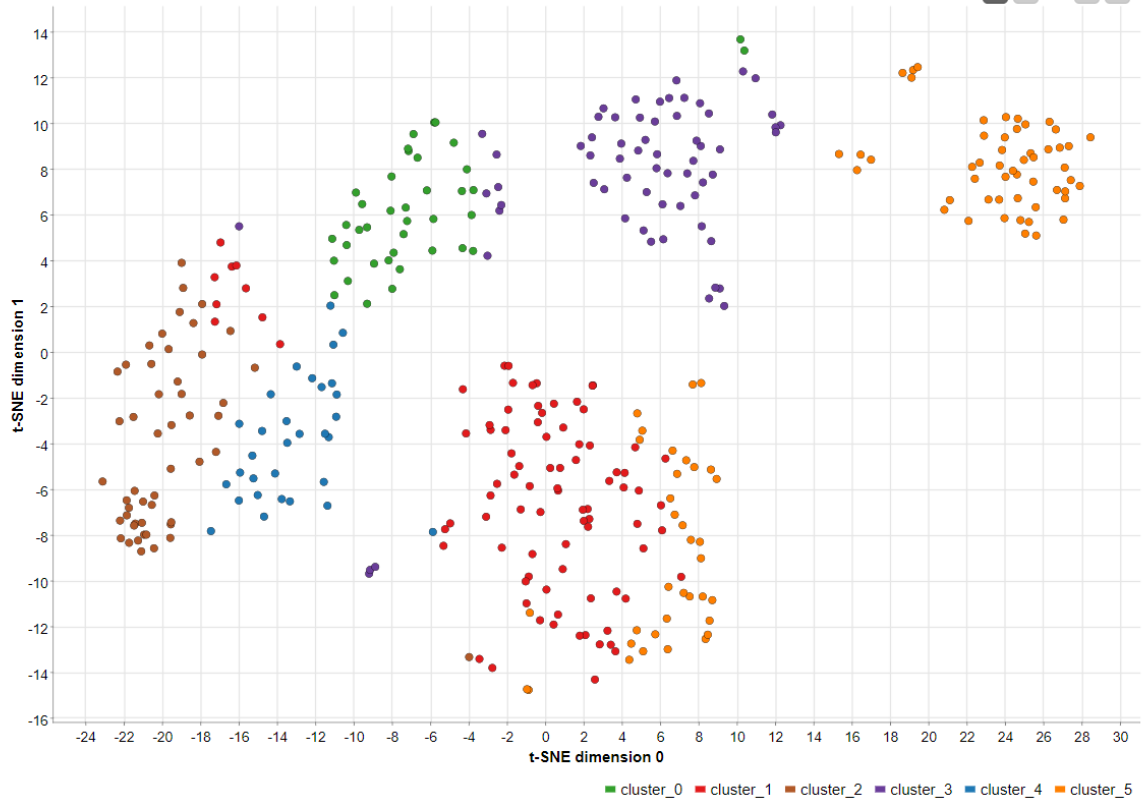
K-means



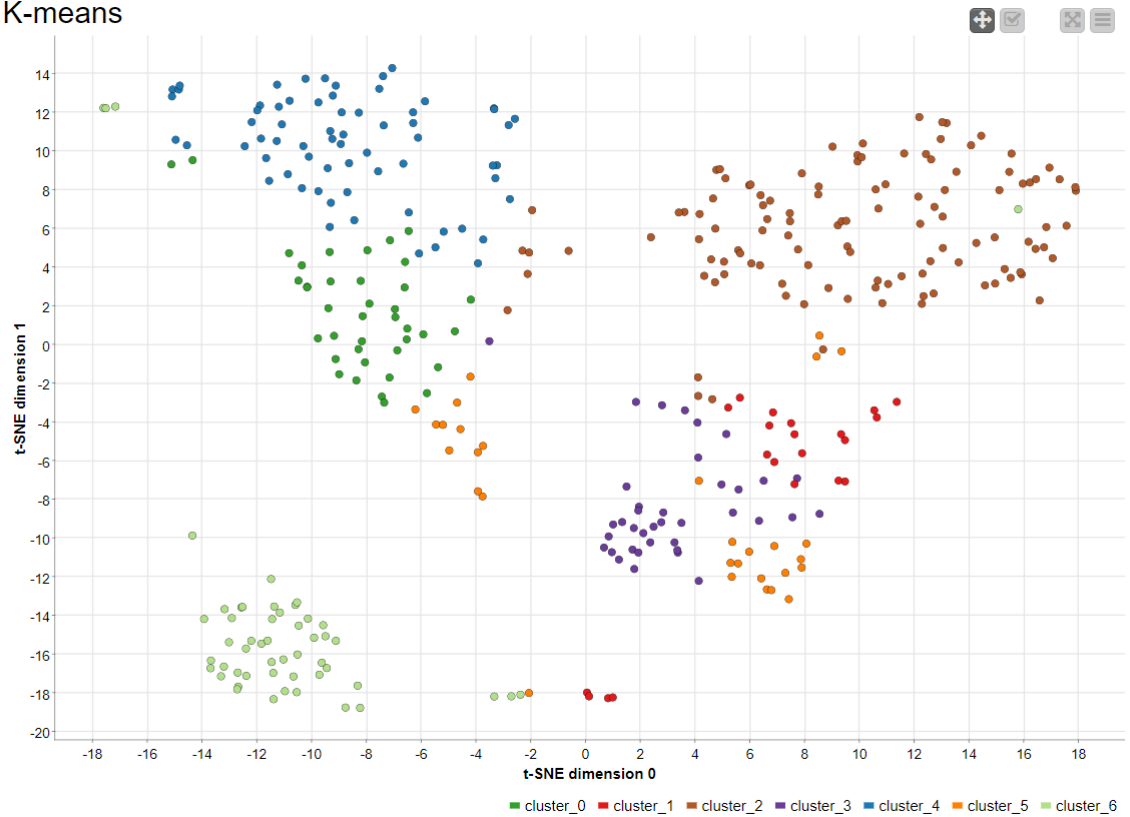
K-means



K-means

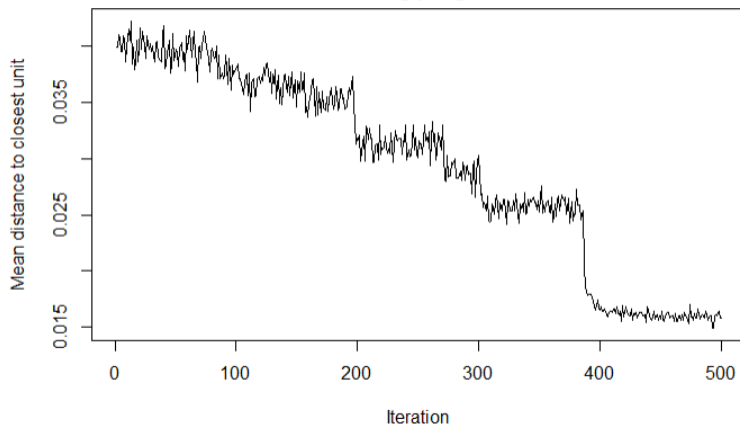


K-means

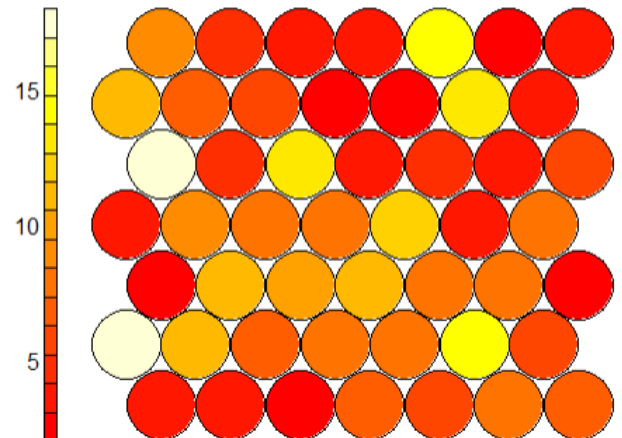


SOM : I choose to work with R as SOM implementation in KNIME do not offer visualization yet. The implementation is done using Kohonen package. I use hexagonal map with 8x8 size and bubble neighborhood function. The learning rate is set between 0.05 and 0.01 with 500 iterations. With this configuration we get 0.015 mean distance to closest node and between 5-10 samples per node .At the end we obtain this result.

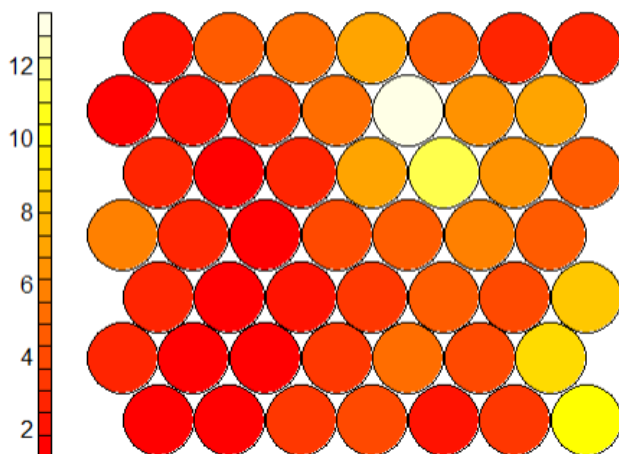
Training progress



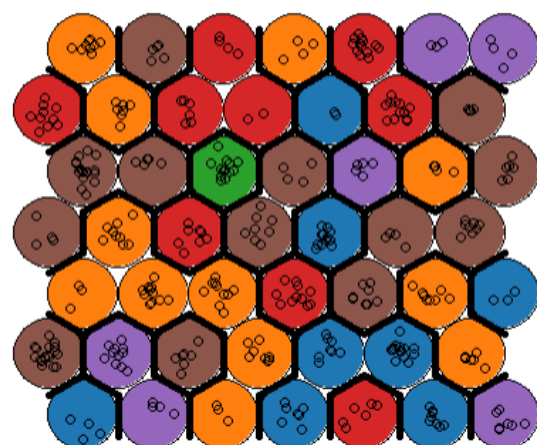
Node Counts



U-matrix



Clusters



Results:

For PCA results we can see that clusters are overlapping (e.g. cluster 2 and 4). Overall, the clusters are not well separated. t-SNE gives good results with clear cluster identifications. For k-means we used several numbers of clusters overall, 5 or 6 clusters suits this dataset. For SOM we can see that cluster melt between each other, and no clear pattern is identified.