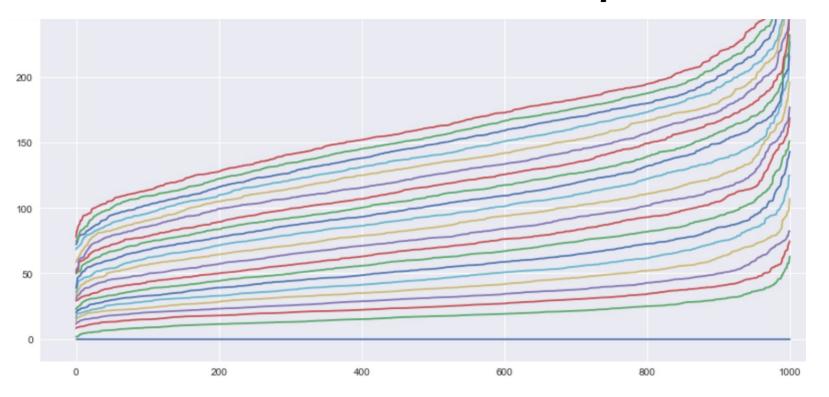# Data visualization & clustering
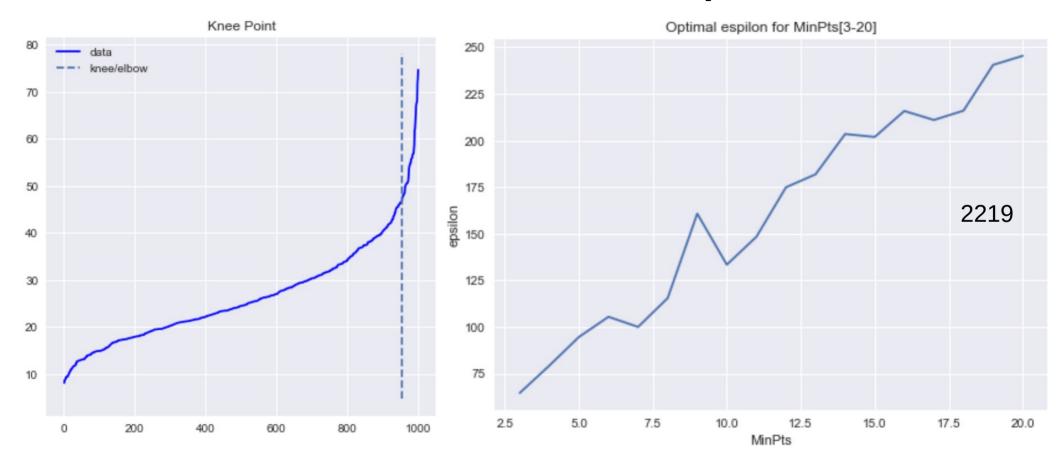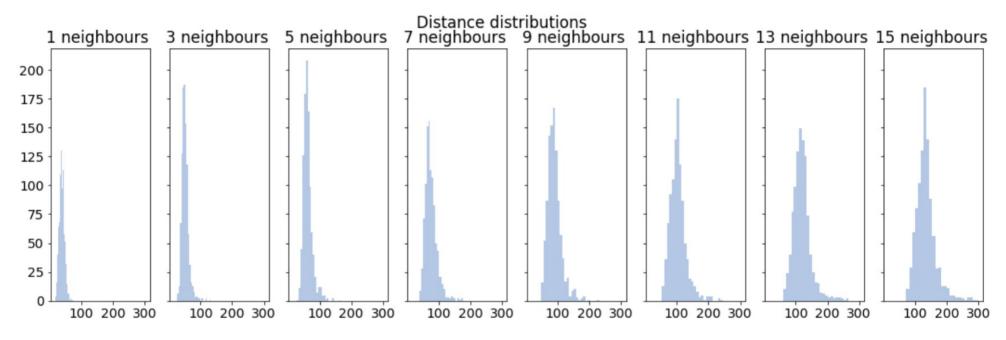
# DBSCAN: find *eps*



2219

To find the threshold for the knee or elbow, the best approach we found was to use the &kneed$ python package. Given x and y arrays, kneed attempts to identify the knee/elbow point of a line fit to the data. The knee/elbow is defined as the point of the line with maximum curvature.
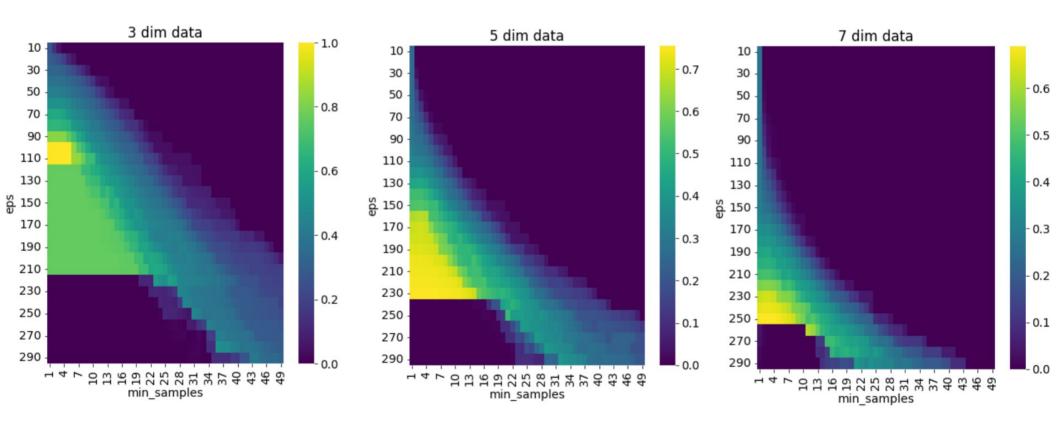
# DBSCAN: find *eps*



2219

# DBSCAN: find *eps*

Distance distributions

Of course, the distance between farer neighbours increases accordingly. Moreover, the distribution seem to have the same shape. As we will see in the following part, DBSCAN performs best with 3 and 5 minpts. And the values of eps choosen will be such that the whole distribution stays on the left of such value.

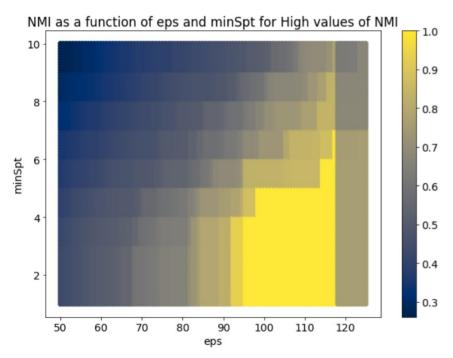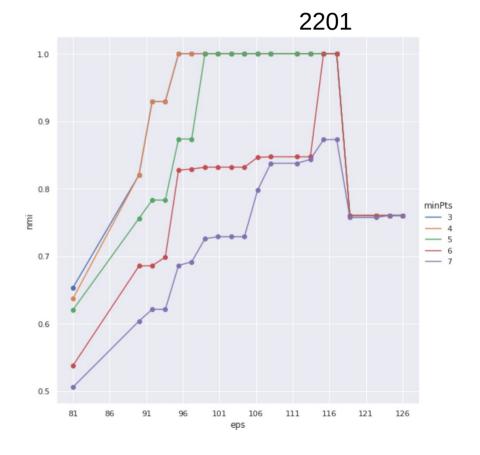# DBSCAN: eps & minPts

2216



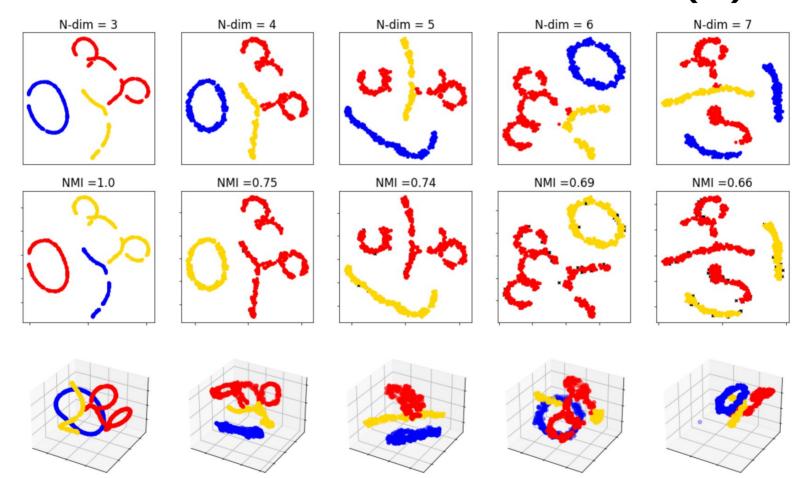Usually, as a rule of thumb, minPts is often set to 2 times the dimensionality of the data.

# DBSCAN: eps & minPts

2211



NMI as a function of eps and minSpt for High values of NMI

2201

# t-SNE and dimensions (L)



2216

# Affiliation Calculation

If two points are strongly affiliated then they are more likly to be in the same cluster. If two points are closer together after T-SNE is applied then they are said to be more strongly affiliated. If they are closer together for many stochasitc itterations of t-SNE then they are also more strongly affiliated. The affiliation of point $p$ between all other points $P$ is calculated as shown below.

$$Affiliation = 1 - norm\left(\sum_i dist(p, P)\right)$$

The distance between two points, $p1 = (x_1, y_1)$ and $P_n = (x_2, y_2)$, is calculated as,

$$dist(p, P_n) = \sqrt{((x_1 - x_2)^2 + (y_1 - y_2)^2)}.$$

dist(p,P) outputs an array of distances between $p$ and every other point $P_n$. These distances are summed for every iteration of t-SNE, $i$. At this point $p$ can be said to be most affiliated with the points where this sum is the lowest. Then they will have been closer together for many iterations of t-SNE. This sum is largly dependant on the number of itterations. To remove this dependancy the result is normalised as shown.

dist(p,P) outputs an array of distances between $p$ and every other point $P_n$. These distances are summed for every iteration of t-SNE, $i$. At this point $p$ can be said to be most affiliated with the points where this sum is the lowest. Then they will have been closer together for many iterations of t-SNE. This sum is largly dependant on the number of itterations. To remove this dependancy the result is normalised as shown.
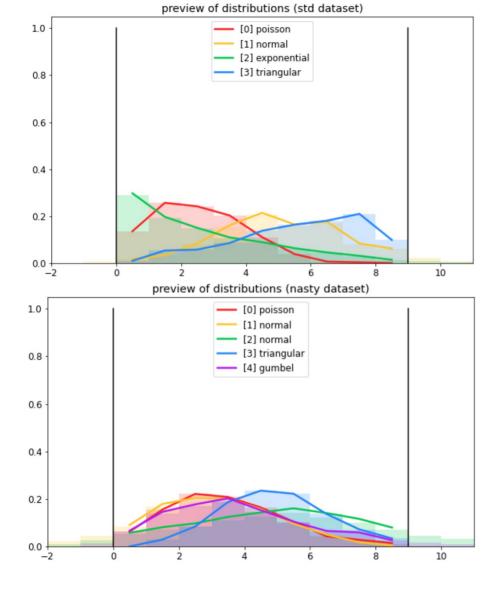
$$norm(data) = \frac{data - min(data)}{range(data)}$$

Where data is an array of numbers. Each element is the sum of the distances of $p$ with $P_n$ over many iterations of t-SNE. After normalisation each element in the array represents the affiliation of $p$ with $P_n$, where numbers close to 0 indicate a strong affiliation and numbers close to 1 indicate a low affiliation. This is reversed by subtracting from 1. This makes affiliation easier to represent graphically.
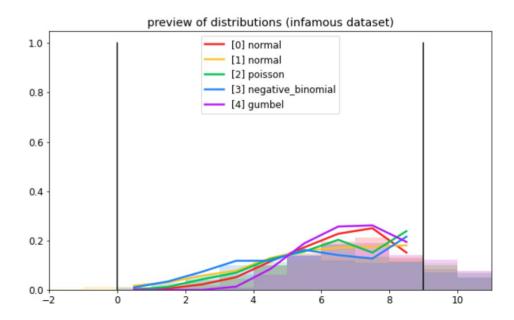
Thus the affiliation of p with another point $P_n$, has been quantised as a number ranging from 0 to 1. Where a value of 1 means they are strongly affiliated and have been placed close together in many iterations of t-SNE and 0 means they are weakly affiliated and have been placed far apart in many iterations of t-SNE.
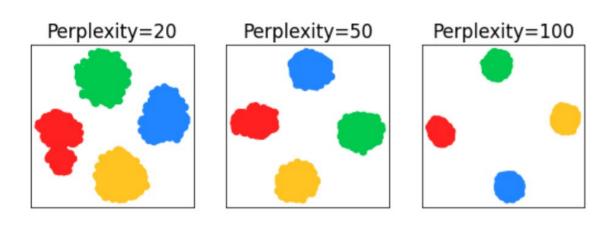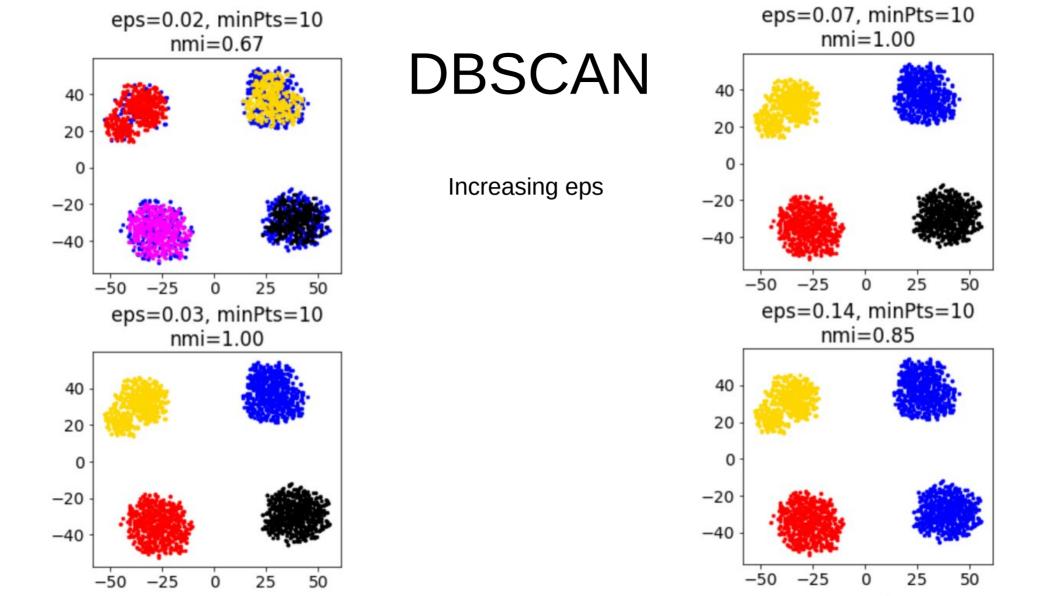
# • Ex.5B: probability distributions

**2202**

# Datasets

preview of distributions (std dataset)

- [0] poisson
- [1] normal
- [2] exponential
- [3] triangular

preview of distributions (nasty dataset)

- [0] poisson
- [1] normal
- [2] normal
- [3] triangular
- [4] gumbel

preview of distributions (infamous dataset)

- [0] normal
- [1] normal
- [2] poisson
- [3] negative_binomial
- [4] gumbel

# T-SNE (Std. dataset)



Perplexity=20 | Perplexity=50 | Perplexity=100

k-means results on t-sne output

```python
import scipy.spatial.distance as spd
```

```python
tsne = manifold.TSNE(n_components=n_components, random_state=123456,
                     perplexity=perplexity, metric=spd.jensenshannon,
                     learning_rate=200, n_iter=1000,
                     n_iter_without_progress=300, min_grad_norm=1e-7,
                     init="random", verbose=1,
                     method="exact", n_jobs=n_jobs, square_distances='legacy')
```

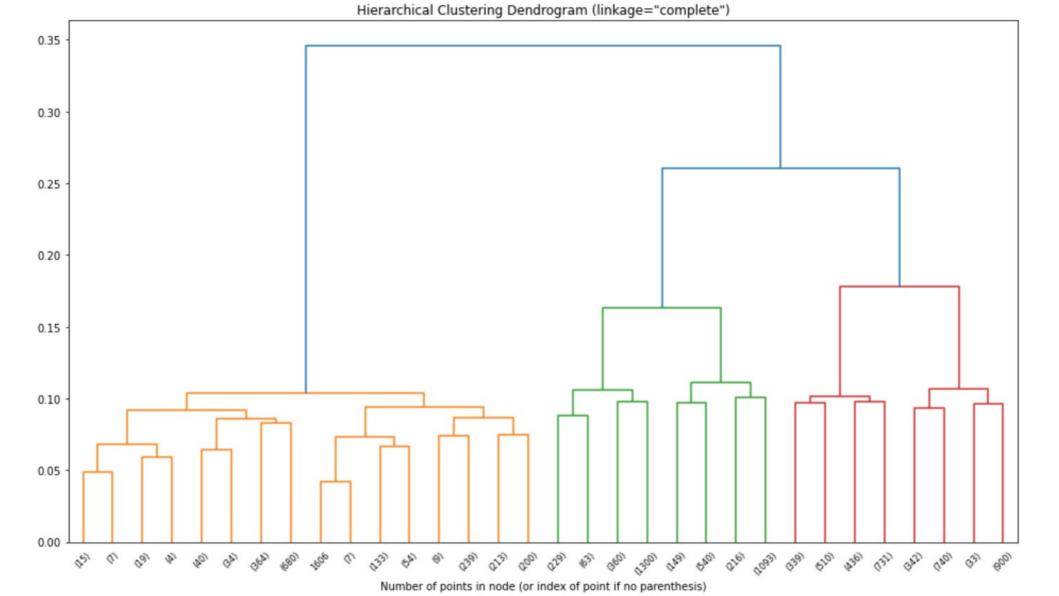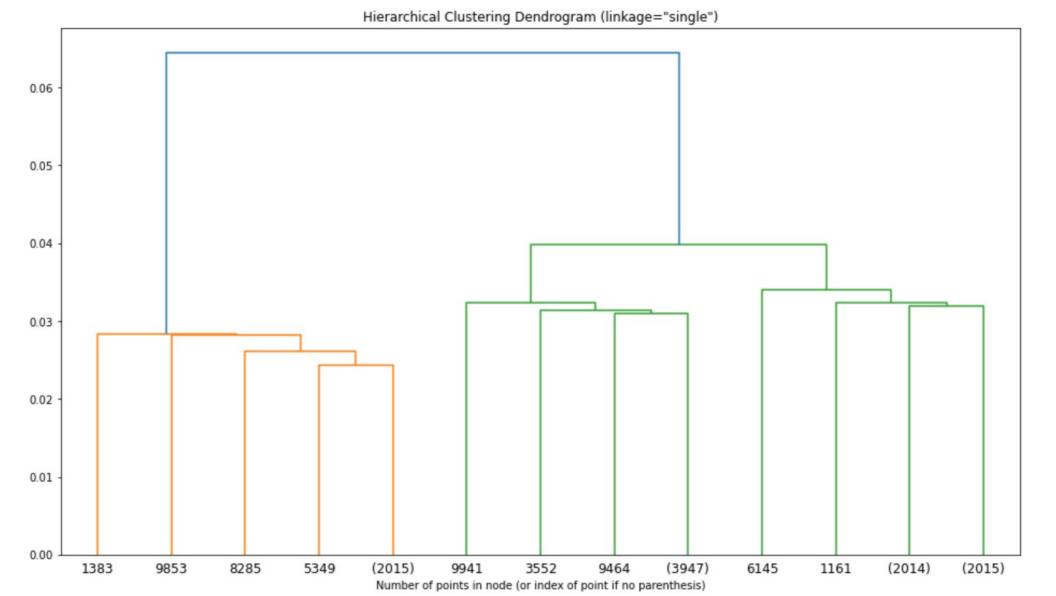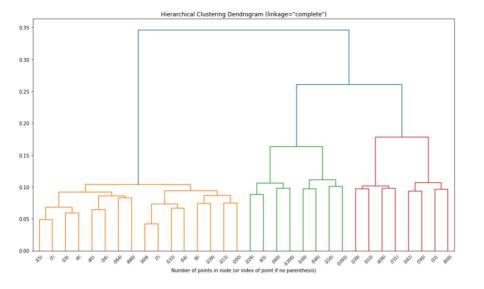DBSCAN

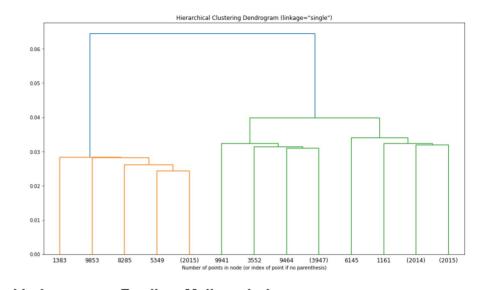Increasing eps

# T-SNE ("infamous" dataset)



From the t-SNE results, we see that the algorithm is able to distinguish all the 5 distributions, even if the dataset has been generated with almost overlapping distributions.
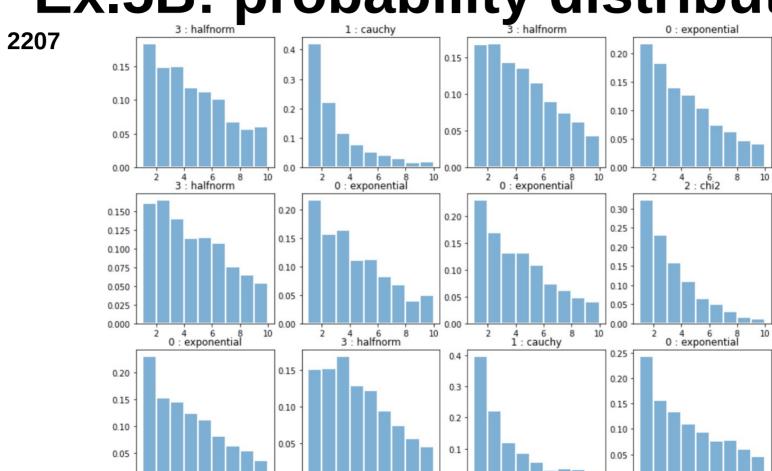
Hierarchical Clustering Dendrogram (linkage="complete")

Hierarchical Clustering Dendrogram (linkage="single")

| Linkage type | Adjusted Rand index score | Fowlkes-Mallows index score |
|:---:|:---:|:---:|
| *complete* | 0.9988 | 0.9990 |
| *average* | 0.9988 | 0.9990 |
| *single* | 0.6166 | 0.7459 |

The results are similar to those obtained on the **nasty** dataset. In particular, the two scoring metrics agree that:
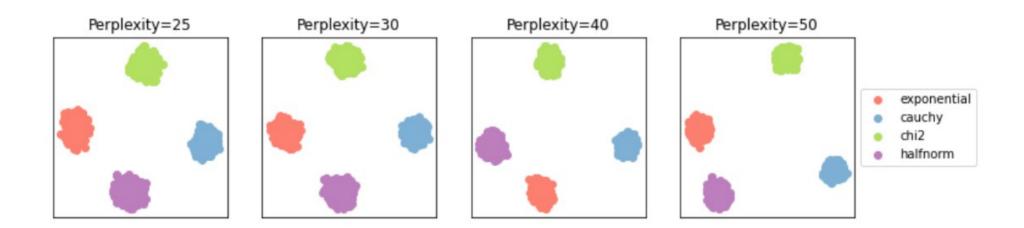
- both *complete* and *average* lead to significantly better results than those obtained with *single*
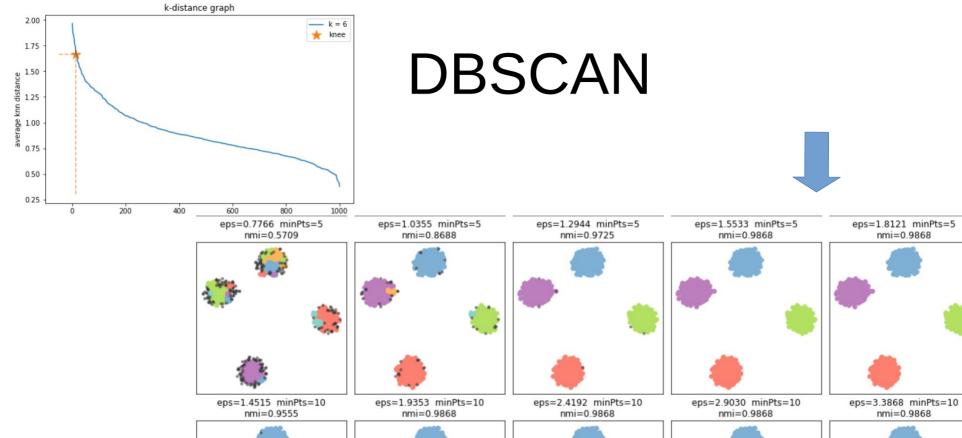- both *complete* and *average* lead to an almost perfect clustering of the **infamous** dataset

# • Ex.5B: probability distributions

**2207**

# t-SNE



Perplexity=25 | Perplexity=30 | Perplexity=40 | Perplexity=50

- exponential
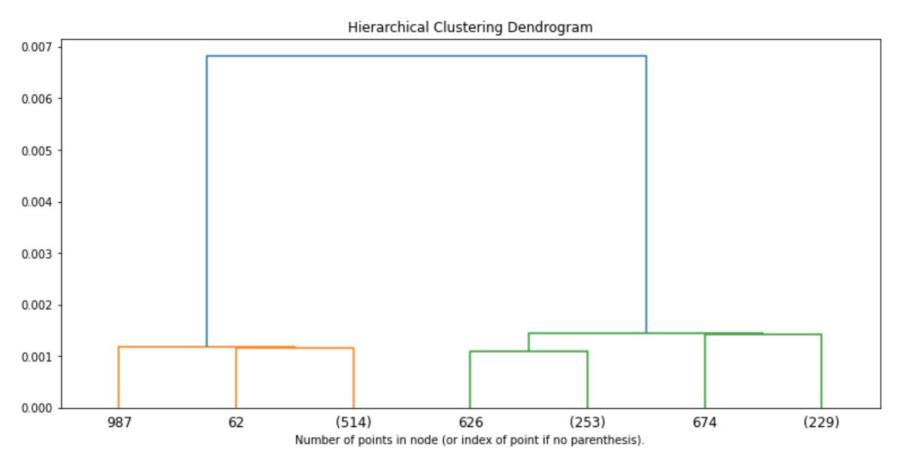- cauchy
- chi2
- halfnorm

# DBSCAN

Hierarchical Clustering Dendrogram

Running the hierarchical clustering over our dataset it is evident what we have noticed from the previous analysis: there are two clusters that are difficult to distinguish with the used metric. As we can notice from the dendogram only three clusters are detected and the situation doesn't change increasing the number of showed levels. So, what we can do is again to use the projected 2D dataset from t-SNE and perfrom the hierarchical clustering using the euclidian metric.

Hierarchical Clustering Dendrogram