

# Clustering and Dimensionality reduction

an informal introduction to  
unsupervised ML

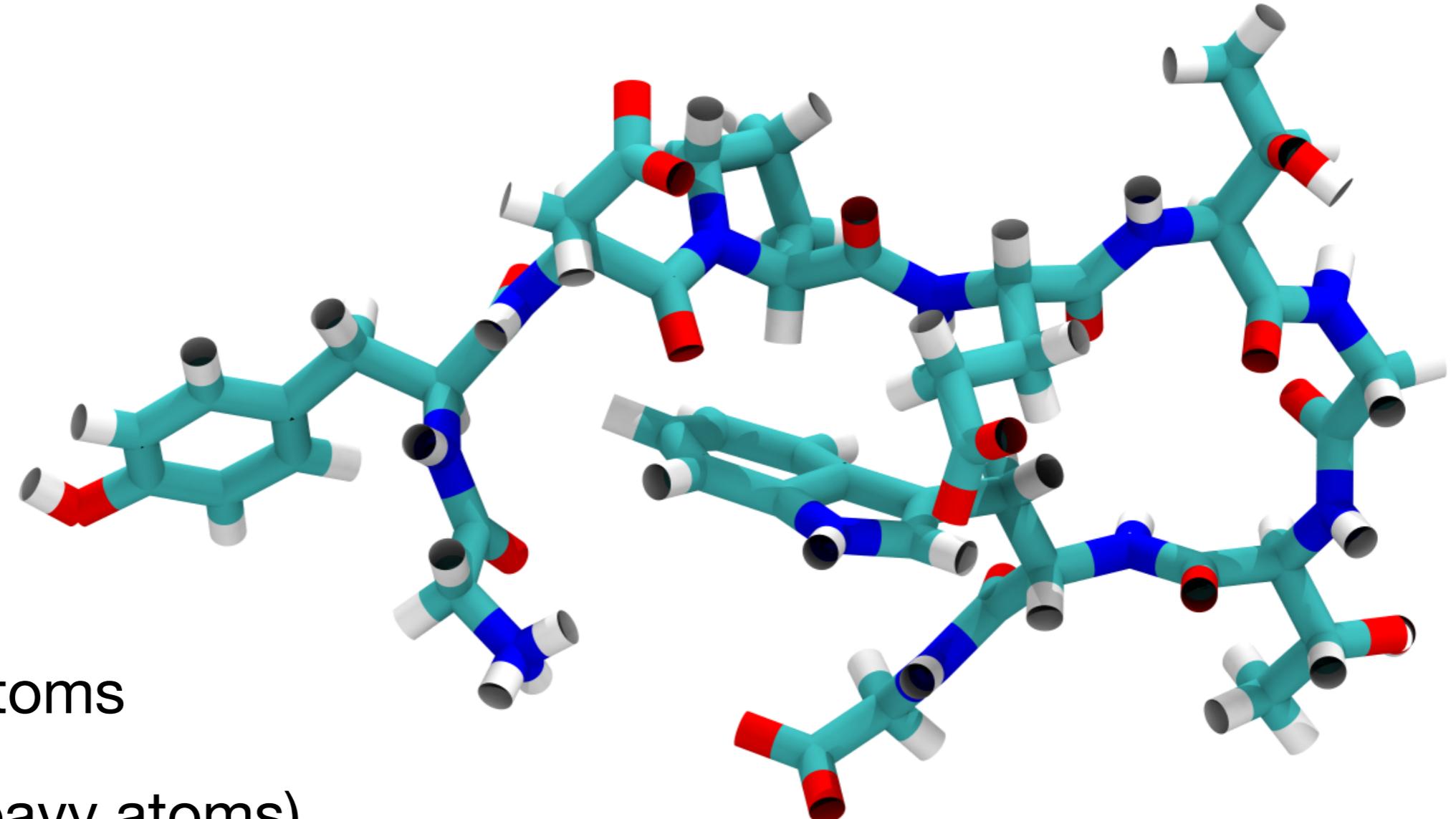
HML

Eugene Klyshko  
Sep 21, 2018

# clustering of protein structures

Protein X:

- 138 atoms
- (83 heavy atoms)
- 10 residues
- 785,612 structures (frames)



# Machine learning

## Supervised

- Classification
- Regression

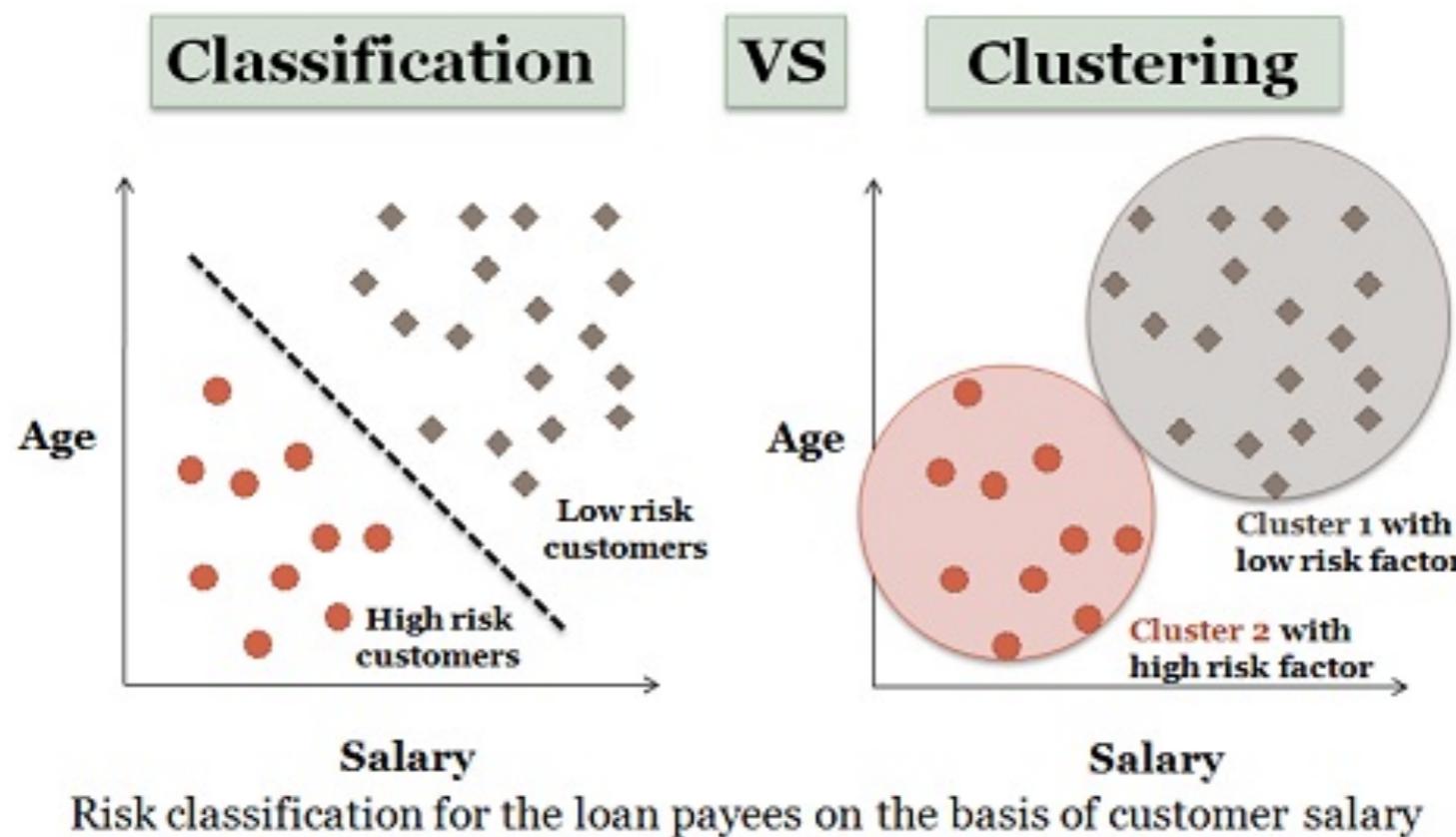
## Unsupervised

- Clustering

**X** - vector of features – **Y** corresponding label

1. given data points: **X** -> **Y**
2. fit and train model
3. determine **Y** for new **X**

1. given data: **X** and no **Y**
2. determine **Y**



Clustering is a technique that performs grouping of data points, so that data is similar within the same cluster, and dissimilar from distinct clusters

# Outline

## 1. Clustering Algorithms:

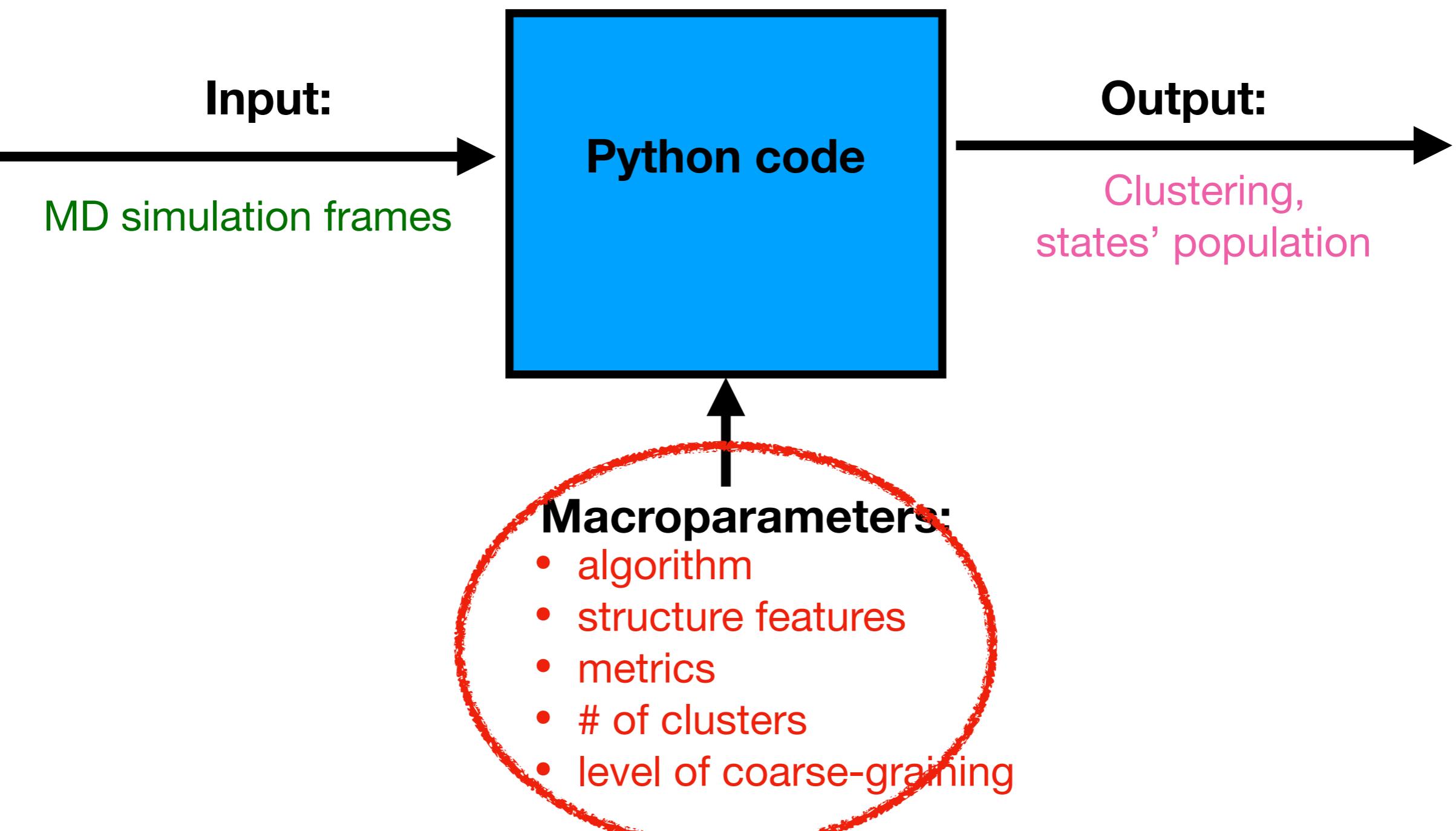
- K-Means
- Hierarchical Clustering

## 2. Dimensionality reduction tools:

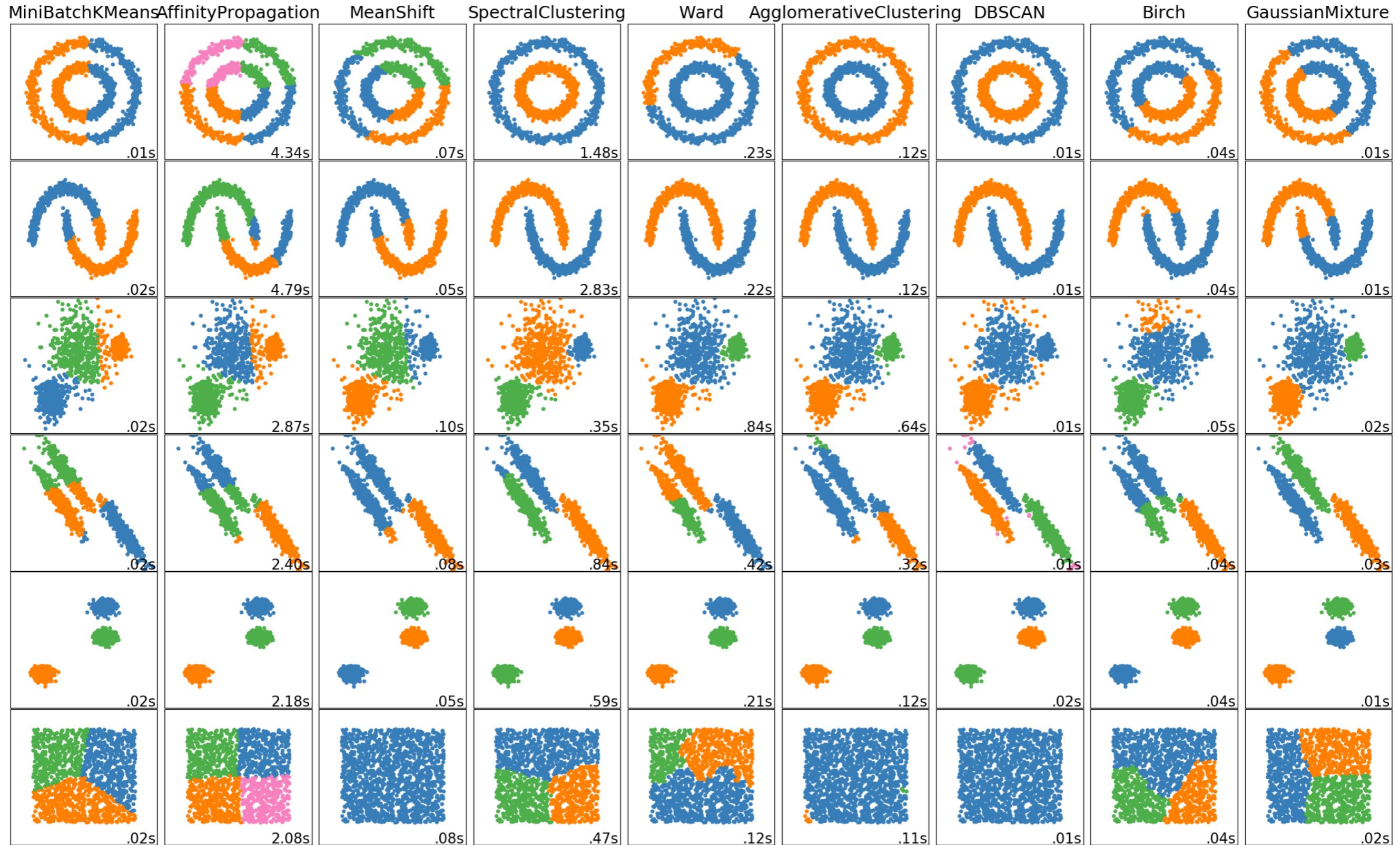
- PCA
- MDS
- tSNE
- UMAP
- Autoencoder NN

# Clustering of structures of unknown protein X

- The main objective is to provide a robust platform for clustering of protein structures from MD trajectories using no prior knowledge



# Clustering algorithms

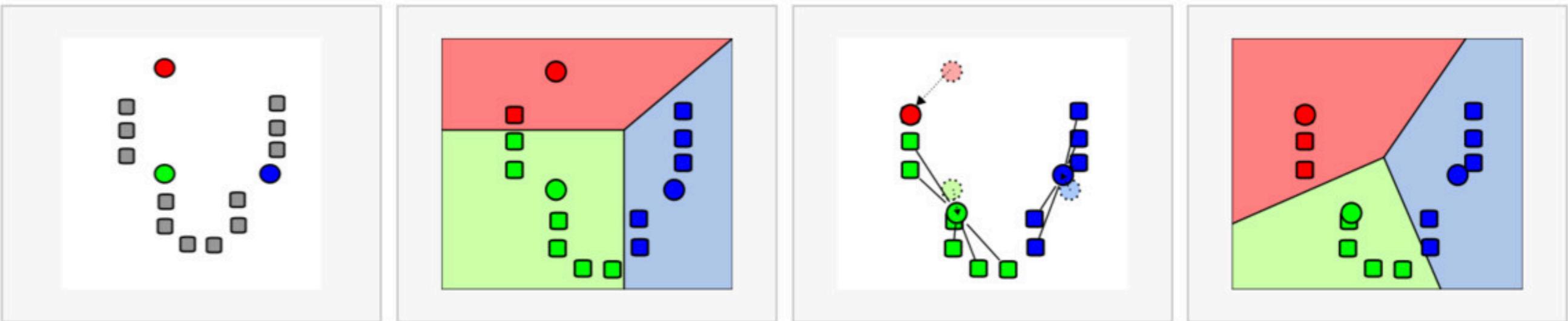


# Macroparameters: algorithm

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large <code>n_samples</code> , medium <code>n_clusters</code> with <a href="#">MiniBatch code</a>	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propagation	damping, sample preference	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium <code>n_samples</code> , small <code>n_clusters</code>	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large <code>n_samples</code> , medium <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes	Distances between nearest points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
Birch	branching factor, threshold, optional global clusterer.	Large <code>n_clusters</code> and <code>n_samples</code>	Large dataset, outlier removal, data reduction.	Euclidean distance between points

# K-Means

Demonstration of the standard algorithm



1.  $k$  initial "means" (in this case  $k=3$ ) are randomly generated within the data domain (shown in color).

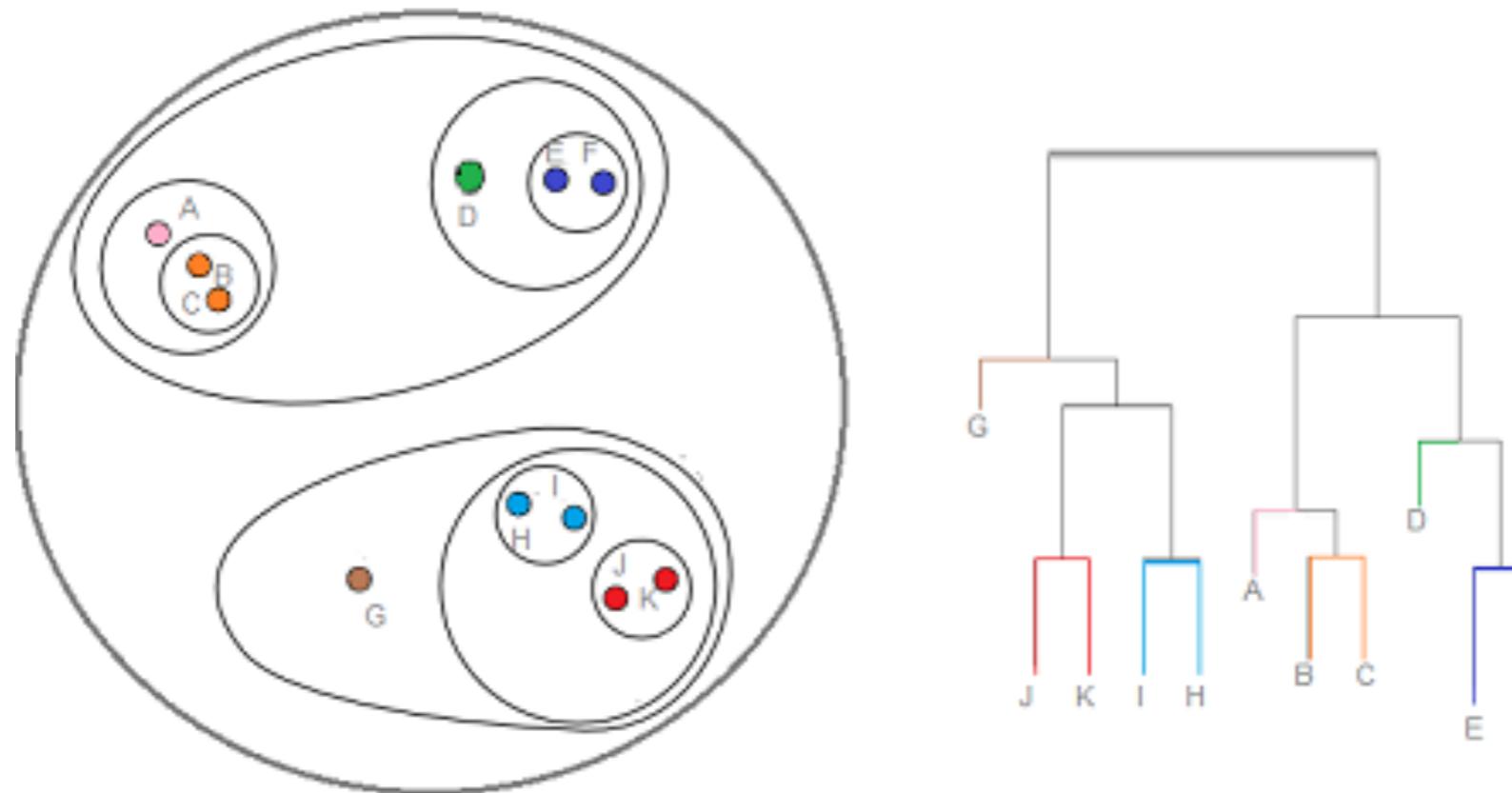
2.  $k$  clusters are created by associating every observation with the nearest mean. The partitions here represent the [Voronoi diagram](#) generated by the means.

3. The [centroid](#) of each of the  $k$  clusters becomes the new mean.

4. Steps 2 and 3 are repeated until convergence has been reached.

# Hierarchical (Agglomerative) Clustering

- We do not know number of clusters, that's why hierarchical clustering algorithm can be a good choice
- Bottom-up approach: each observation starts in its own cluster, and then pairs are merged until we get 1 cluster containing all points



- The standard algorithm for hierarchical agglomerative clustering has a time complexity of  $O(n^3)$  and requires  $O(n^2)$  memory, which makes it too slow for even medium data sets.

# Hierarchical (Agglomerative) Clustering

- Metric is how we calculate distance between feature vectors in high-dimensional space

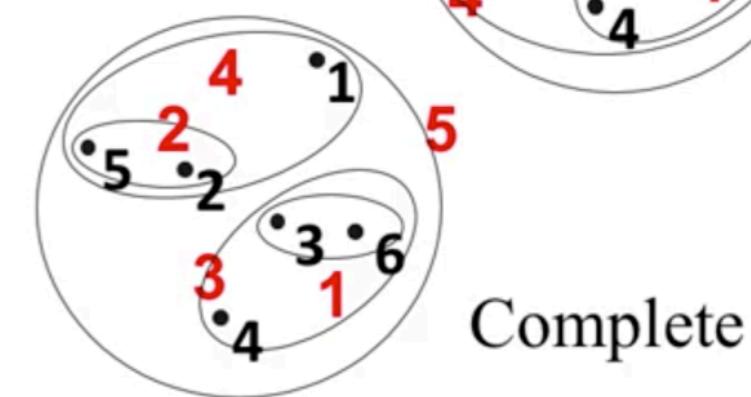
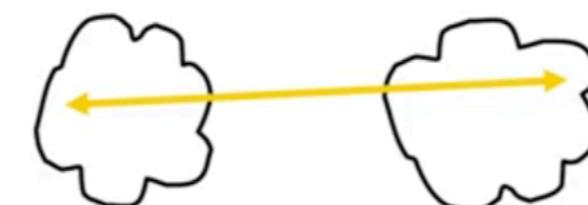
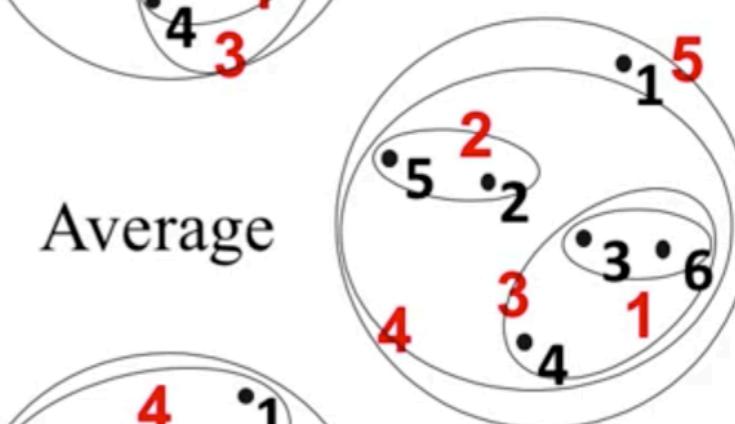
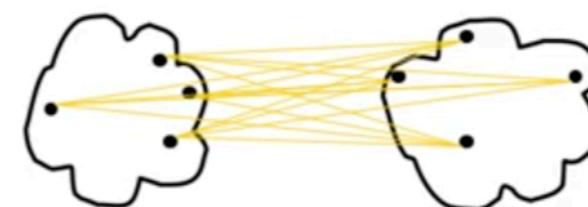
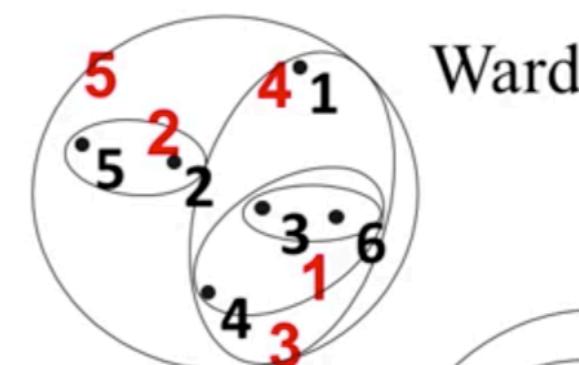
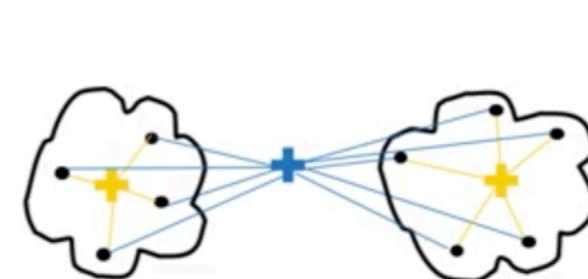
Names	Formula
Euclidean distance	$\ a - b\ _2 = \sqrt{\sum_i (a_i - b_i)^2}$
Squared Euclidean distance	$\ a - b\ _2^2 = \sum_i (a_i - b_i)^2$
Manhattan distance	$\ a - b\ _1 = \sum_i  a_i - b_i $
Maximum distance	$\ a - b\ _\infty = \max_i  a_i - b_i $
Mahalanobis distance	$\sqrt{(a - b)^\top S^{-1} (a - b)}$ where $S$ is the Covariance matrix

# Hierarchical (Agglomerative) Clustering

- Linkage is the method used to measure distance between newly created clusters

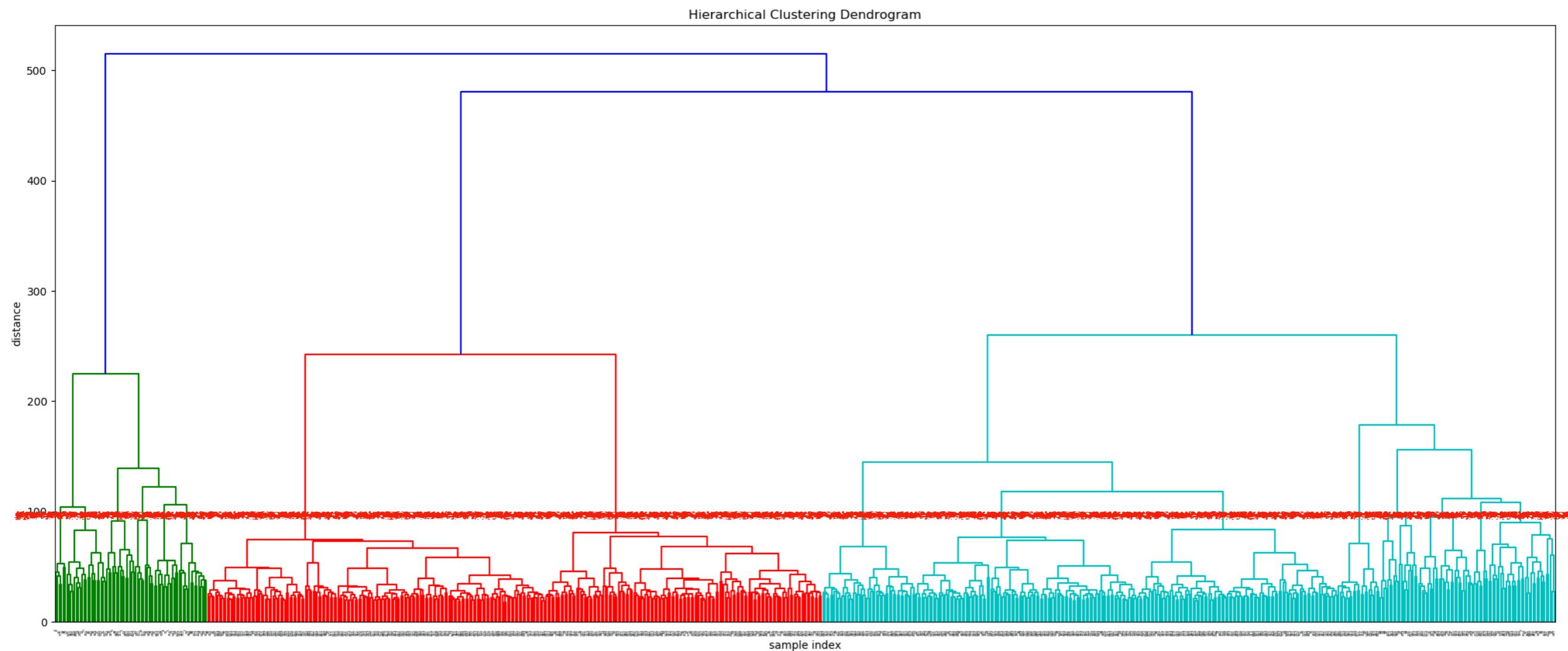
- **Ward's method**

- Least increase in total variance (around cluster centroids)

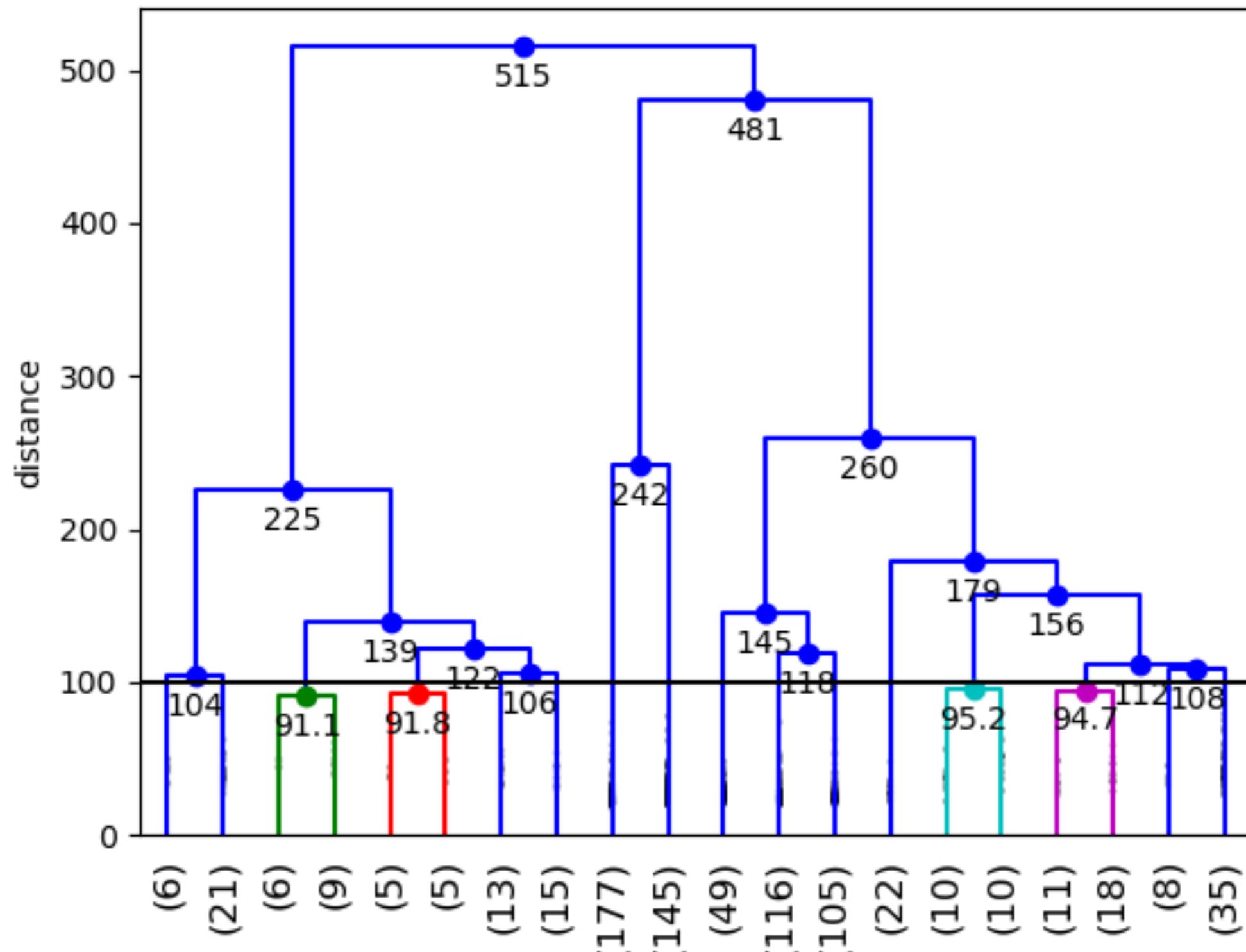


\*cophenetic index to compare efficiency of linkage criteria

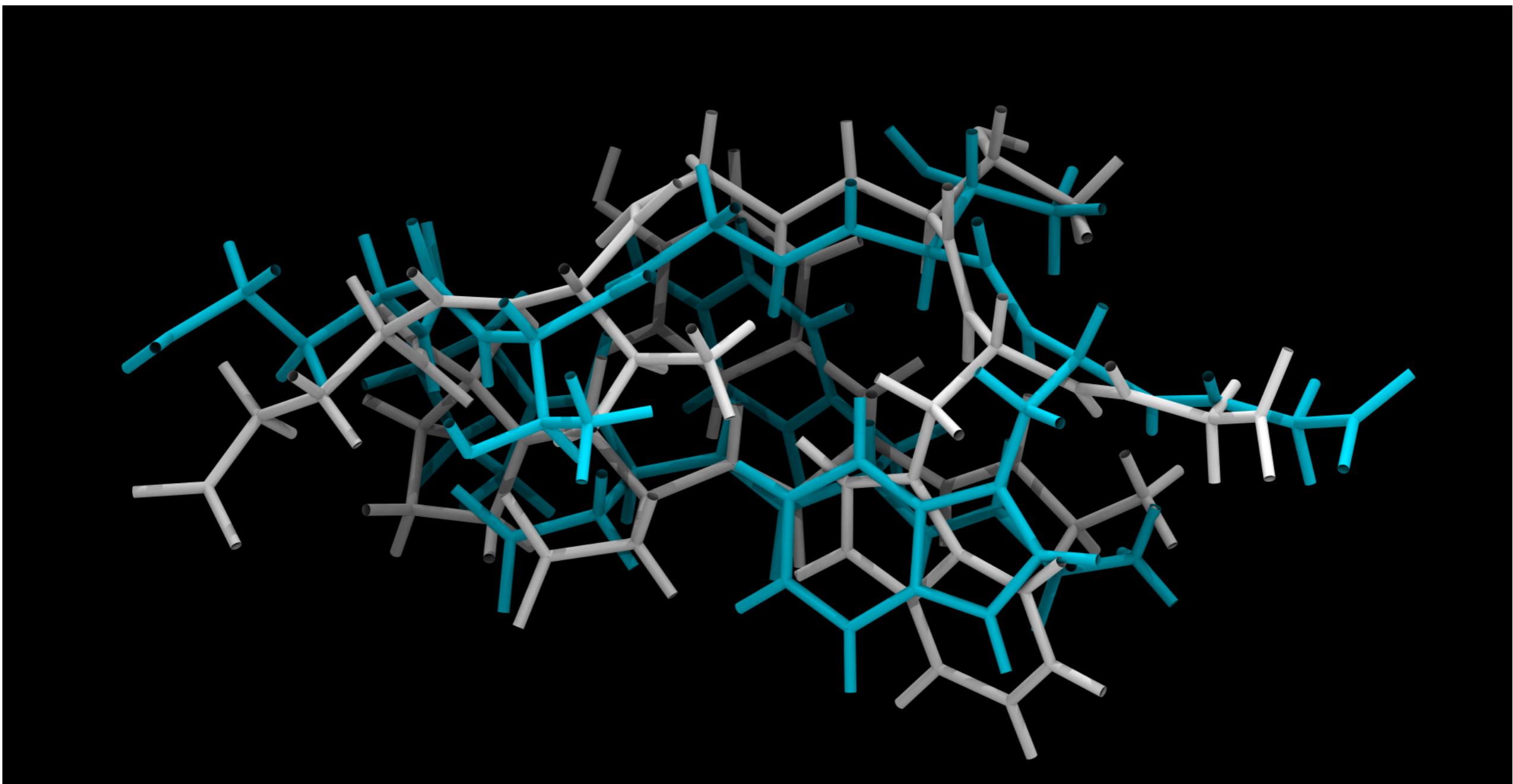
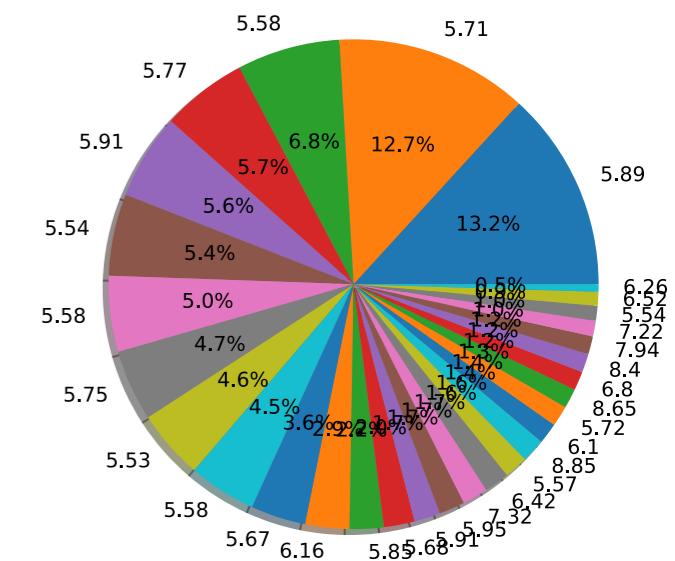
# Full dendrogram



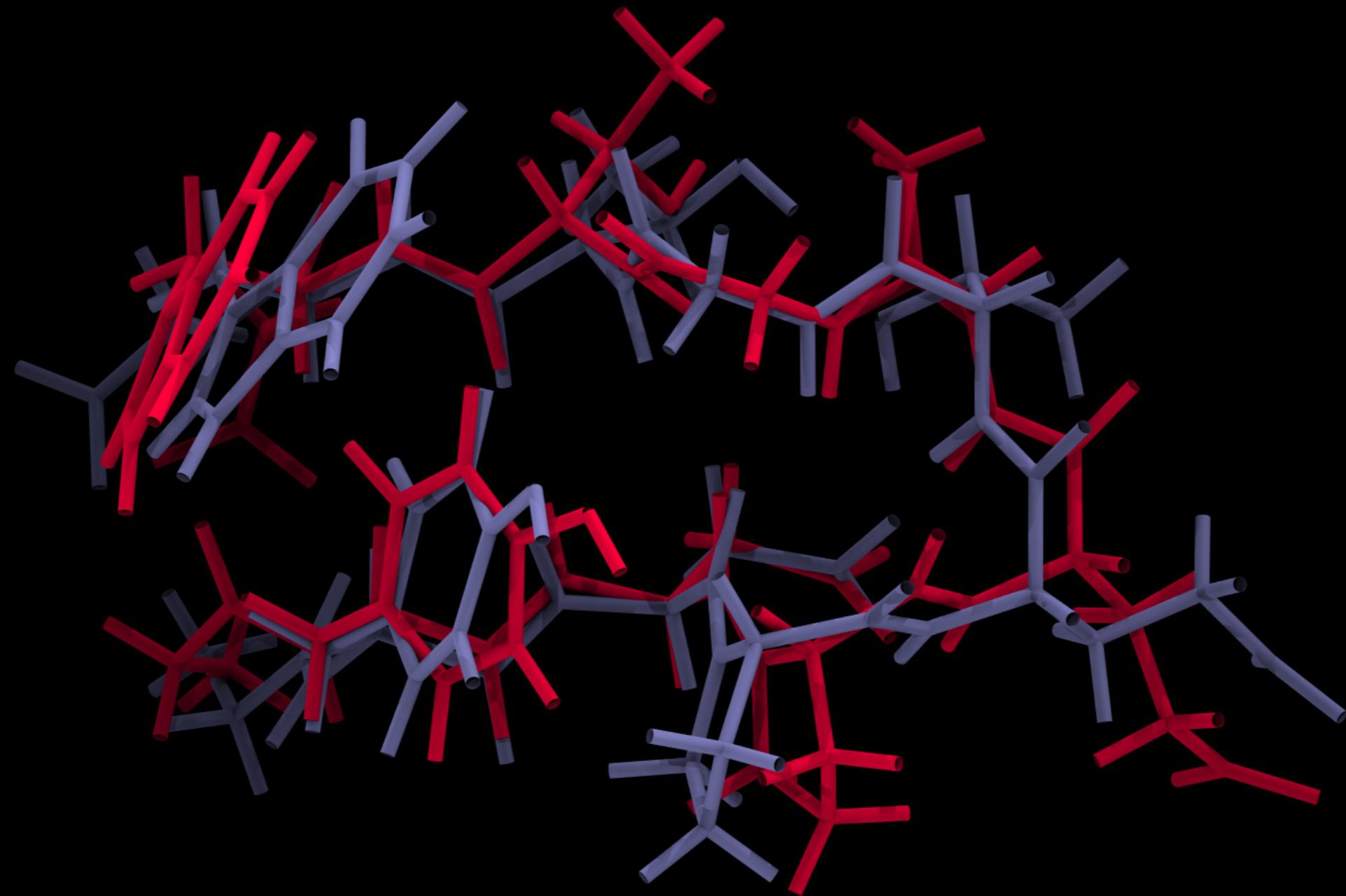
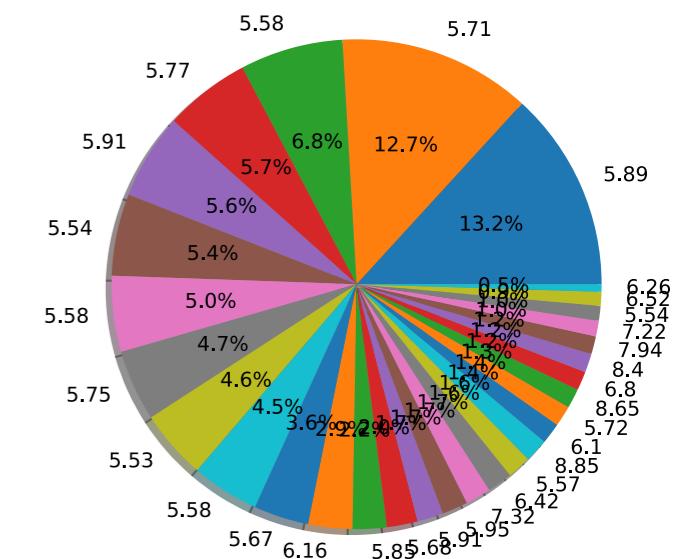
### Hierarchical Clustering Dendrogram (truncated)



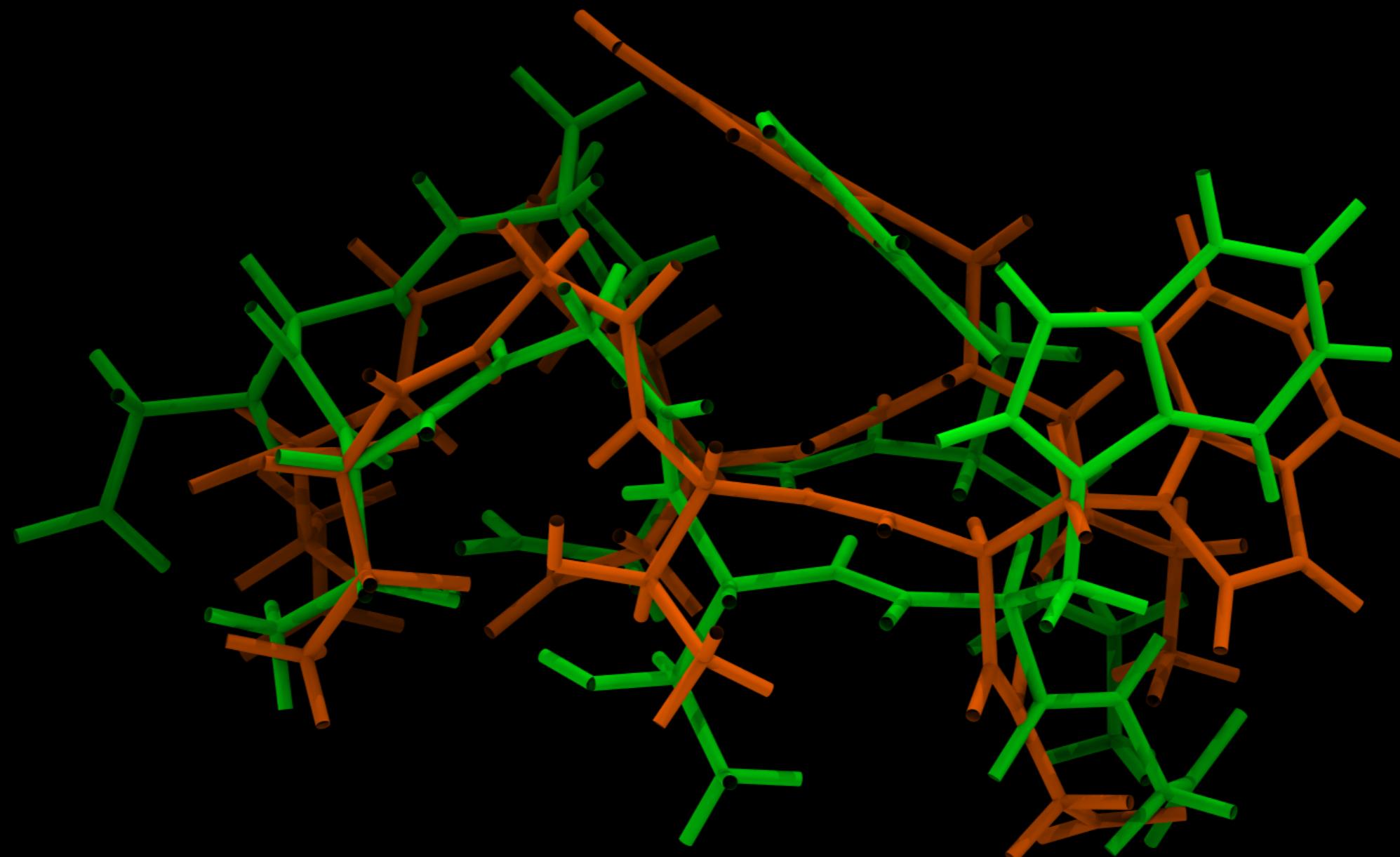
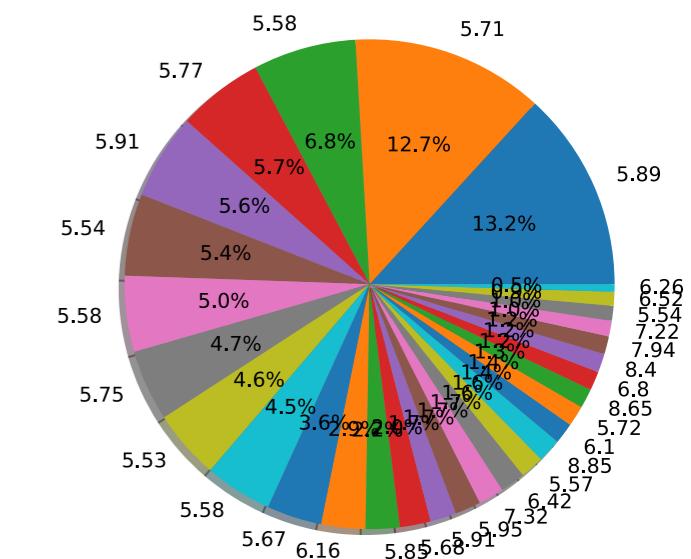
# Visualizing (1-2)



# Visualizing (3–4)

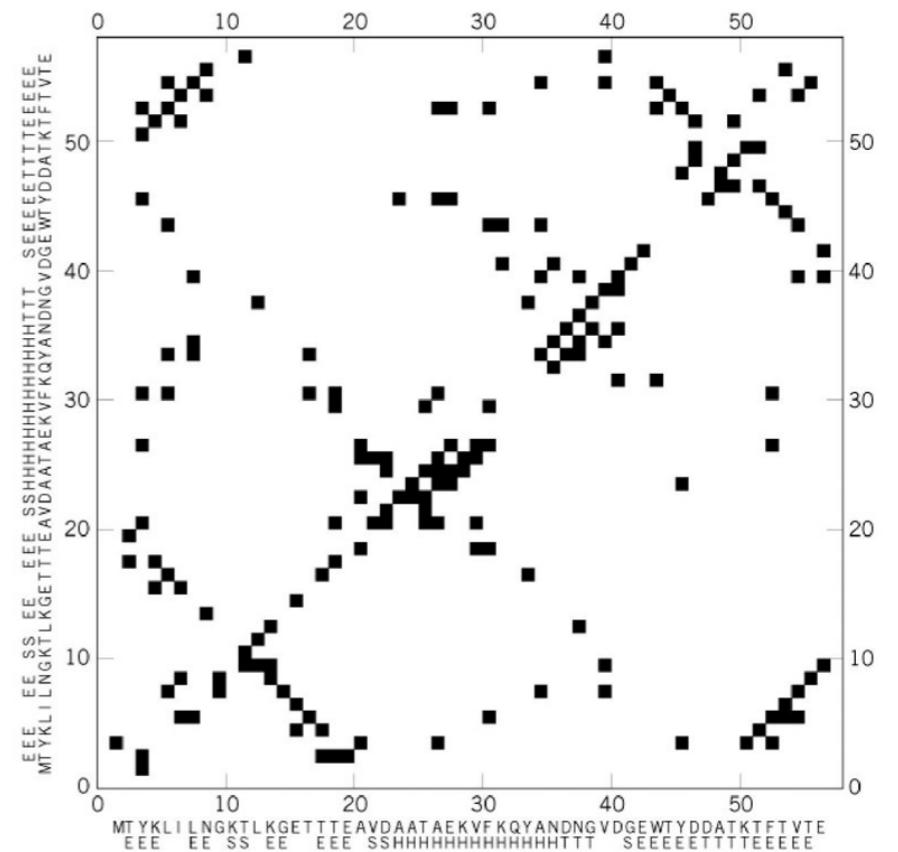


# Visualizing (5–6)



# Feature vector

- Last time we agreed that contact map vector should be changed. Instead of determining the contact for every heavy atom, now we consider contact for each amino acid only (except itself and two neighbouring)
  - If any two heavy atoms in the amino acid are closer than cutoff  $5.5 \text{ \AA}$ , then **1**, else **0**. Total size 36 for  $n = 10$  amino acids.
  - (Let's also add radius of gyration as the 37th feature)



# Feature vector - contact map

Total dimension = 36  
~~12,000~~

# of data points = 780,000

# of unique data points = 16,000

$\mathbf{X} = \begin{vmatrix} 1 \\ 0 \\ \dots \\ 0 \\ 1 \end{vmatrix}$

- I used ‘cosine’ metrics - I’ll explain why
- the result doesn’t change much in terms of linkage ‘method’
- The results will be shown for ‘ward’ linkage

# Metrics: cosine

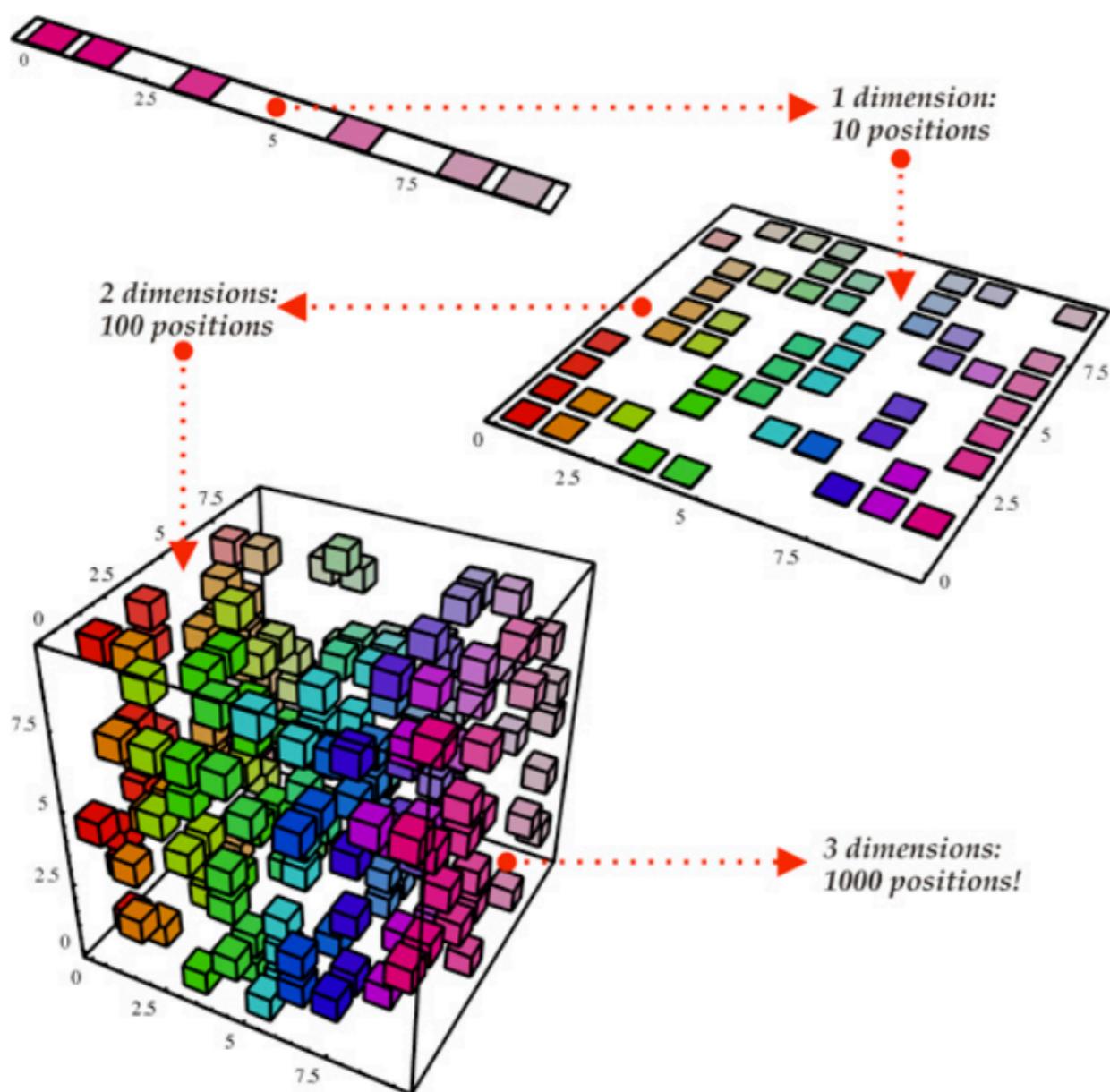
- Euclidean distance has the property that true values and false values for features carry the same weight, and this does not work well for contact maps, which are generally sparse.
- When dealing with sparse data sets in high dimensions, it is often the case that the presence of an attribute is more significant than its absence.

the **cosine** distance between vectors  $u$  and  $v$ :  $d(u, v) = 1 - \frac{u \cdot v}{\|u\|_2 \|v\|_2}$

**It is thus a judgment of orientation and not magnitude**

# Dimensionality reduction

- We often need this to represent multidimensional data in lower dimensional space:

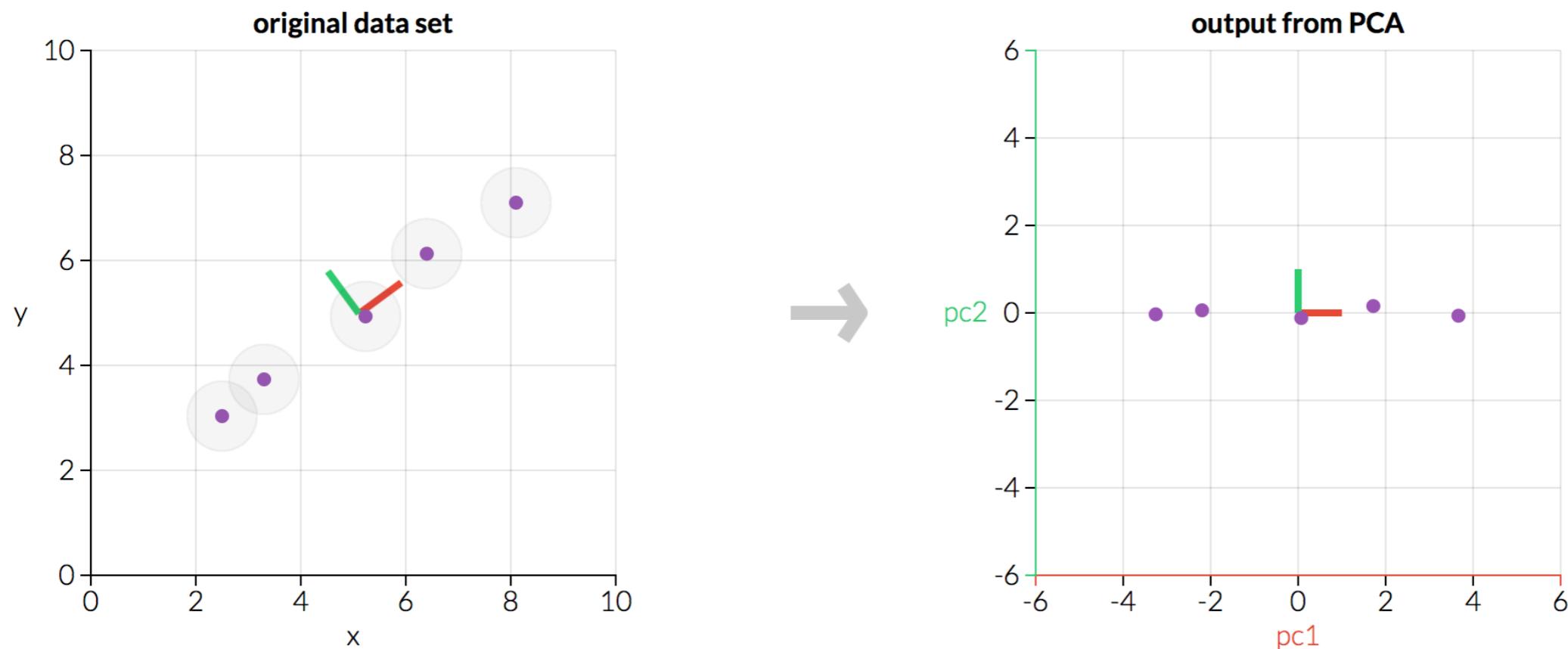


36D → 2D

## Why do we need it?

- visualization
- speed up computations
- representation
- simplification
- get insights about the data

# Dimensionality reduction: PCA - Principal Component analysis

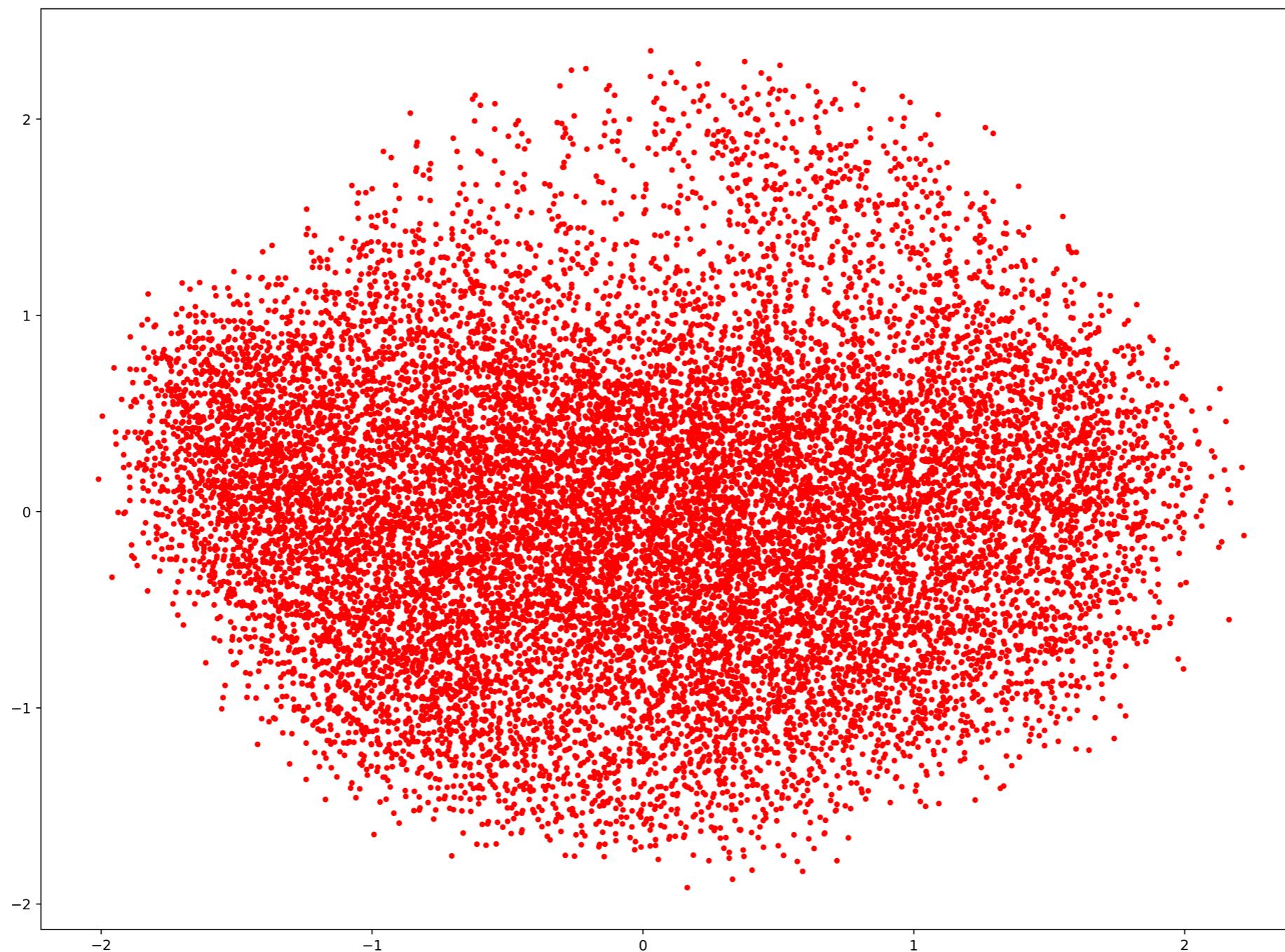


- Finding the directions of the greatest variance of data: i.e. find eigenvalues of covariance matrix, sort them in descending order
- Find projections of the data on the eigenvectors, corresponding to the largest k eigenvalues
- $\mathbf{W}_L$  is a matrix made up of the eigenvectors of correlation matrix

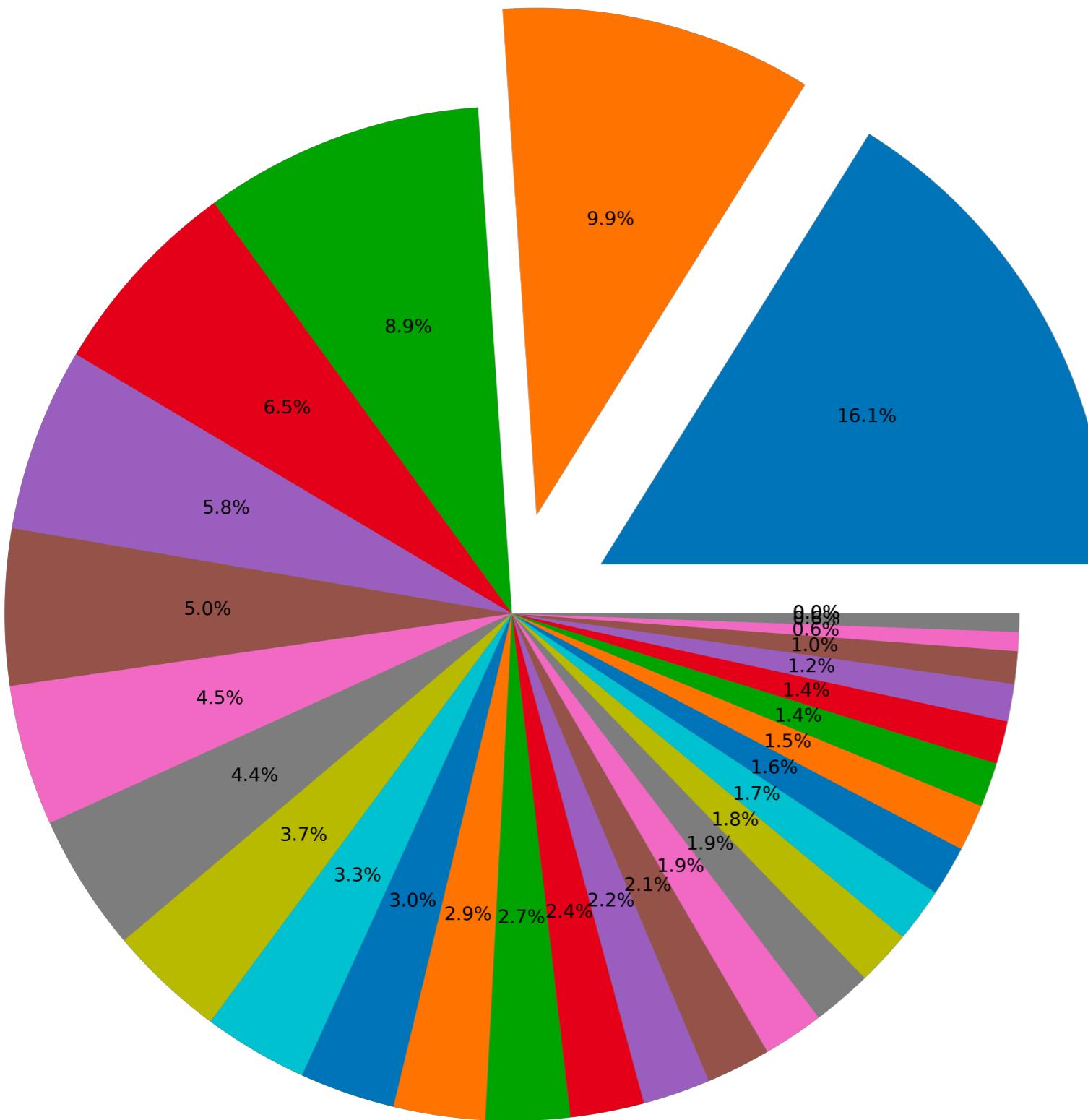
$$\mathbf{T}_L = \mathbf{X}\mathbf{W}_L$$

# PCA - continue

- Our original data (16,000 points) represented by first 2 PCAs:

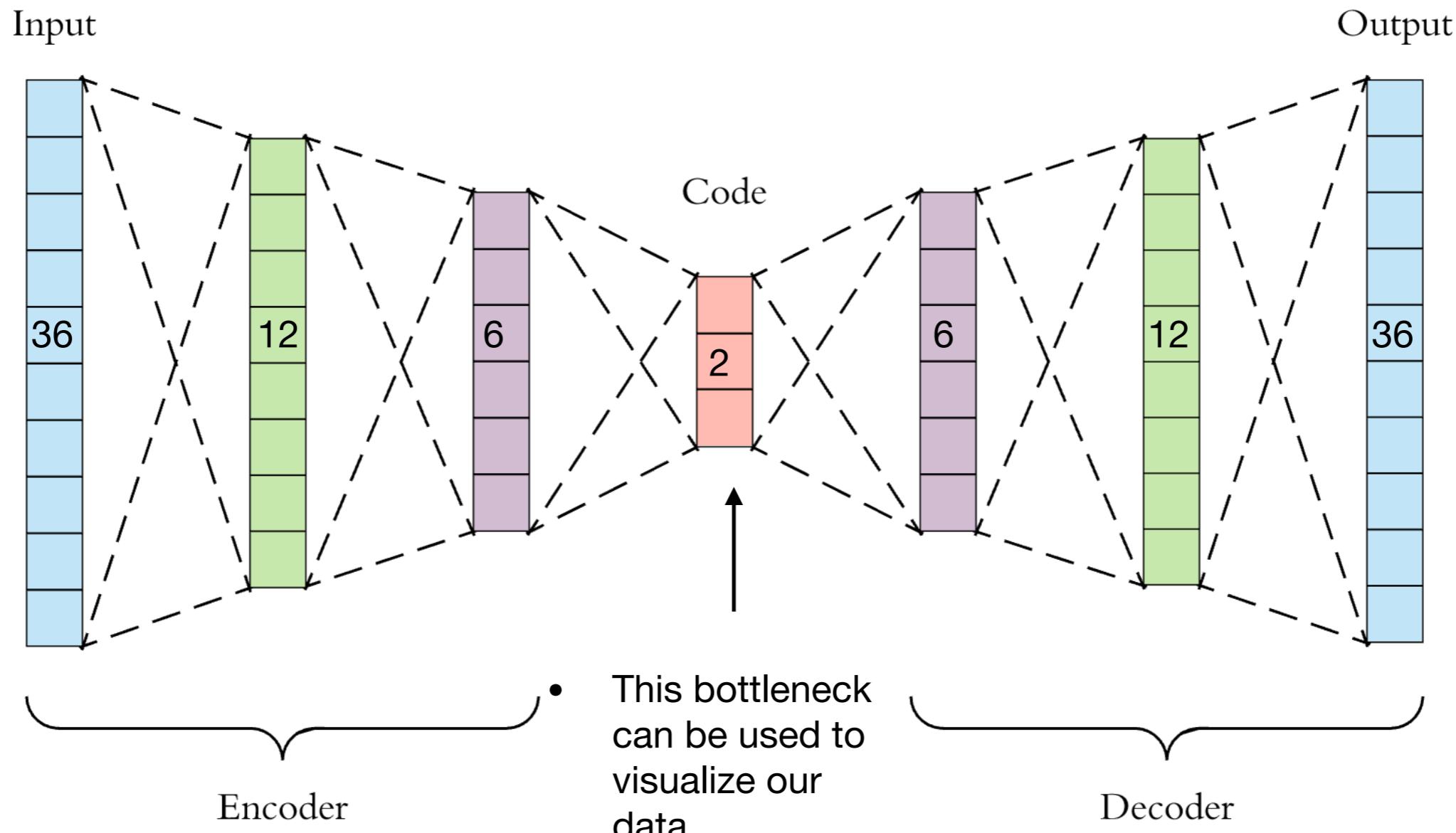


# PCA - continue even more



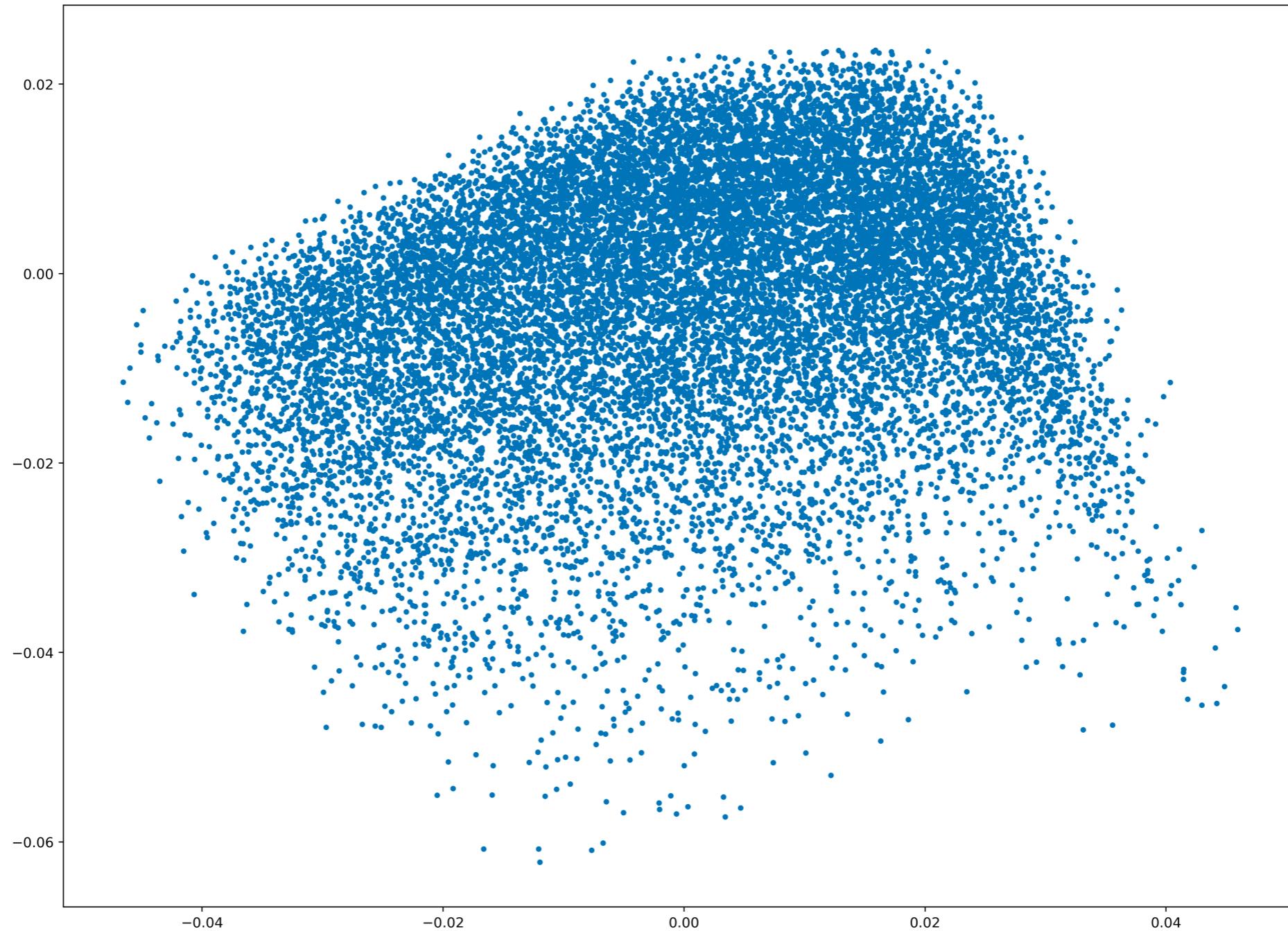
- As you see, 2 first principal components make up only 26% of total information of 36 dimensions
- Hence, we need some other approach
- NN!

# Dimensionality reduction - autoencoder Neural Network



- Used Keras with TensorFlow to build NN: tanh activation functions and regularization are used
- Trained on our 16,000 data points

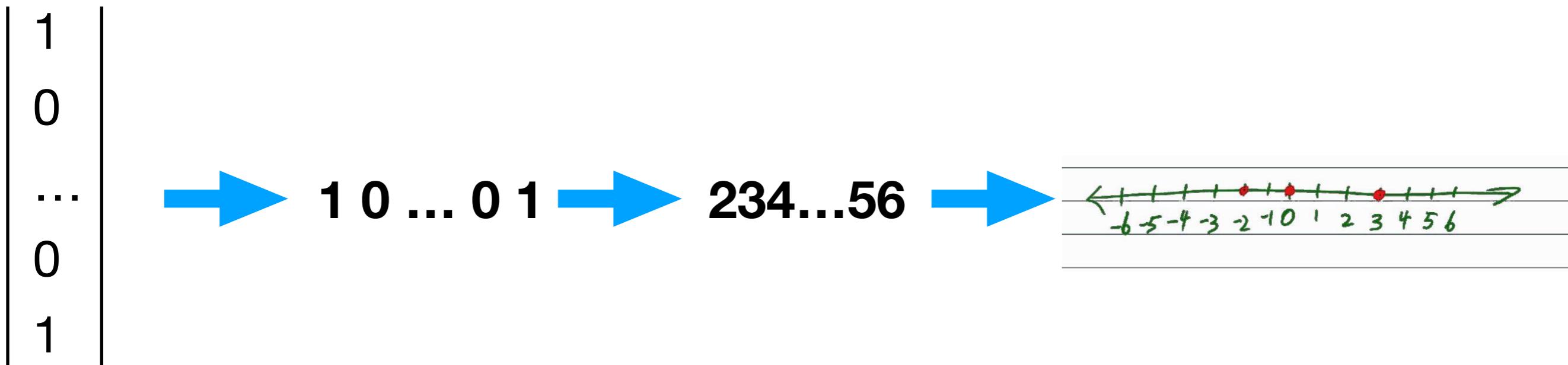
# Autoencoder Neural Network - continue



- 85% accuracy (in other words, it's 85% of information)
- Unfortunately, it doesn't give us any ideas about the number of clusters

# My (maybe) stupid idea

- Our 36D vector is not quite in 36 dimensions of real numbers, right?
- Features can be only 0 or 1
- So let's represent every data point as a binary 36-bit number.
- Computers can deal with 64-bits
- Convert to 64-bit decimal integer
- Plot all integers on the coordinate axis. Voila! 1D is enough



The problem is that we need not only visualization, but the information about the pairwise distances. It should conserve

# Dimensionality reduction - SNE

Stochastic Neighbour Embedding (SNE)

$p_{j|i}$  is a conditional probability for point  $x_i$  to pick  $x_j$  as its neighbour [original multidimensional space]

$q_{j|i}$  is a conditional probability for  $y_i$  to pick  $y_j$  as the neighbour [lower dimensional space]

Cost function [KL-divergence]

Gradient of the cost function for minimization

Gradient decsend method

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2/2\sigma_i^2)},$$

$$q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}.$$

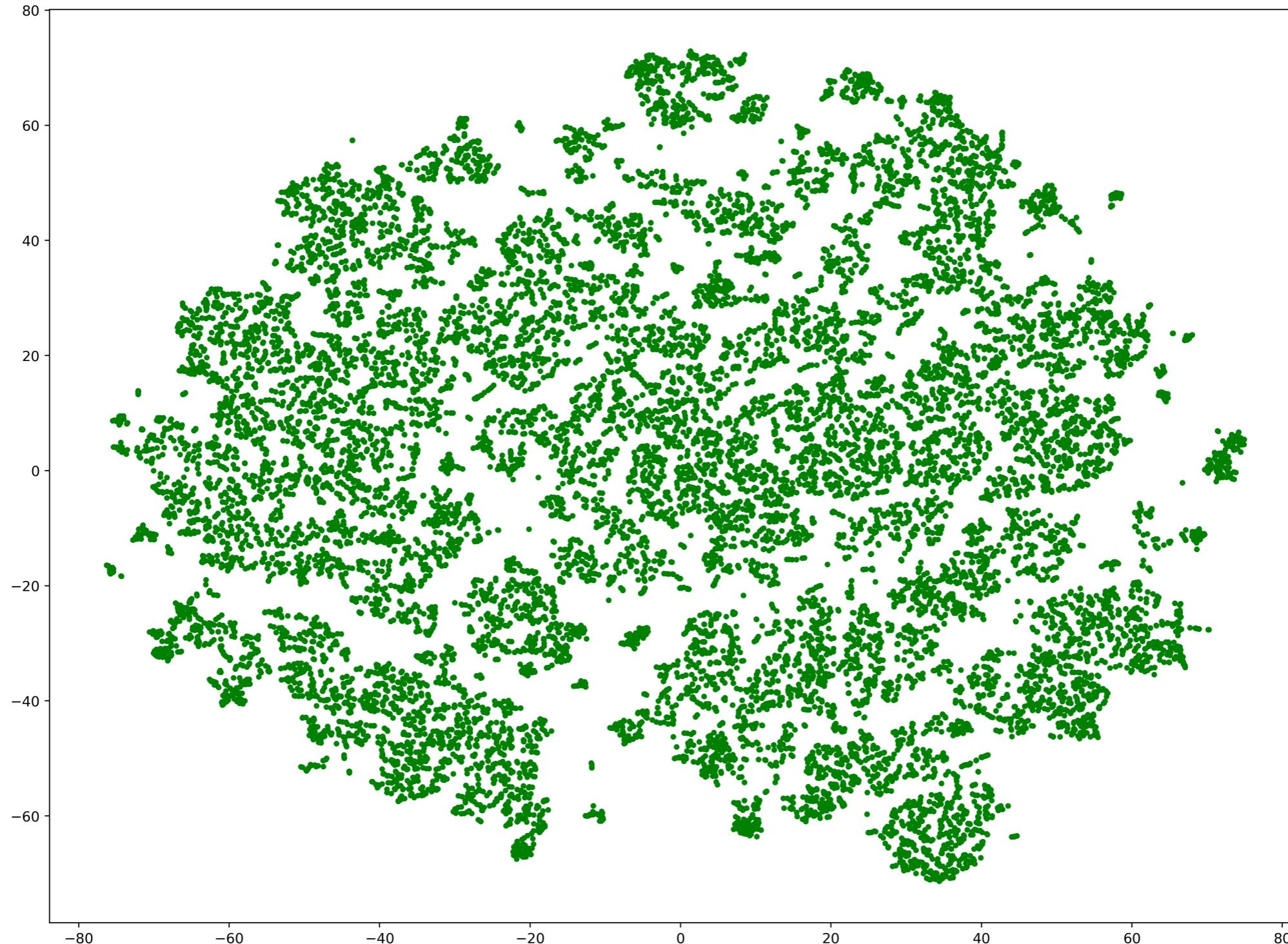
$$C = \sum_i KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}},$$

$$\frac{\delta C}{\delta y_i} = 2 \sum_j (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})(y_i - y_j).$$

$$\gamma^{(t)} = \gamma^{(t-1)} + \eta \frac{\delta C}{\delta \gamma} + \alpha(t) (\gamma^{(t-1)} - \gamma^{(t-2)})$$

t-SNE - modified SNE

# tSNE - continue



Great tool for visualization and determining number of clusters