



CLUSTERING

MOHAMMAD GHODDOSI

WHAT IS CLUSTERING

- grouping a set of objects
- objects in the same group are more similar to each other than to those in other groups.

CLUSTERING VS CLASSIFICATION

- Clustering is unsupervised
 - We just have input
 - We don't know possible outputs
- Classification is supervised
 - We have input and label
 - We know possible outputs

K-MEANS

Input:

$D = \{t_1, t_2, \dots, t_n\}$ // Set of elements

K // Number of desired clusters

Output:

K // Set of clusters

K-Means algorithm:

Assign initial values for m_1, m_2, \dots, m_k

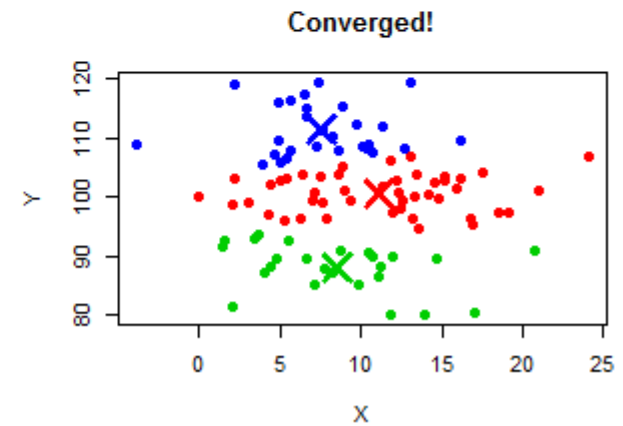
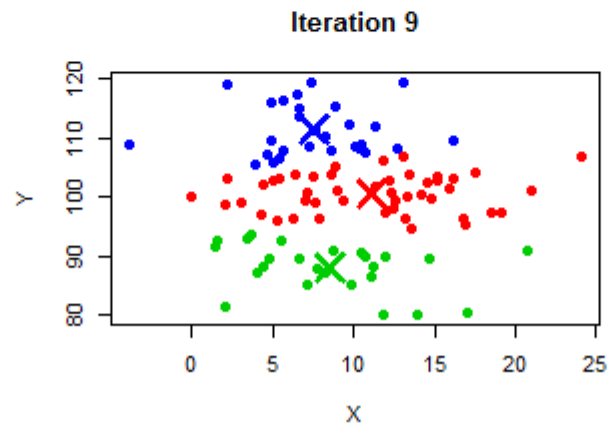
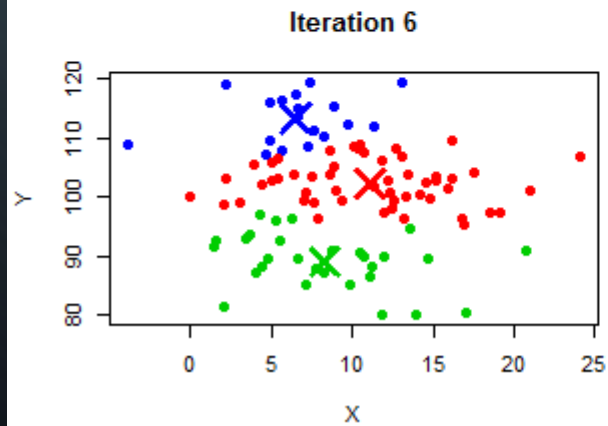
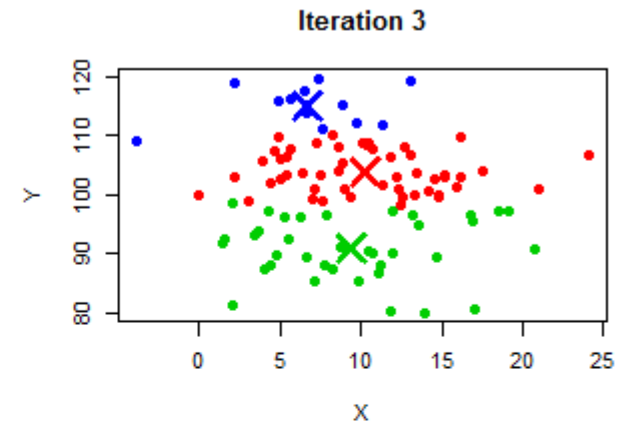
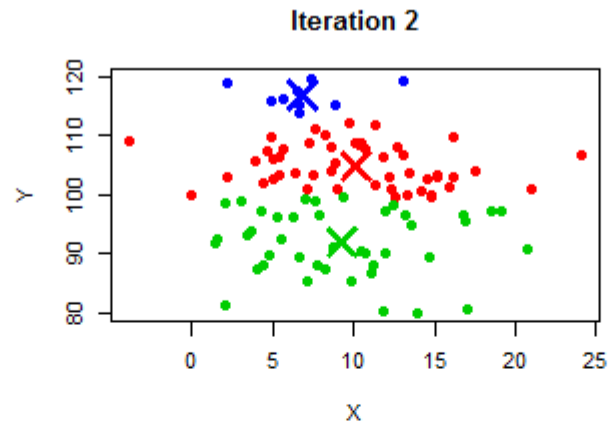
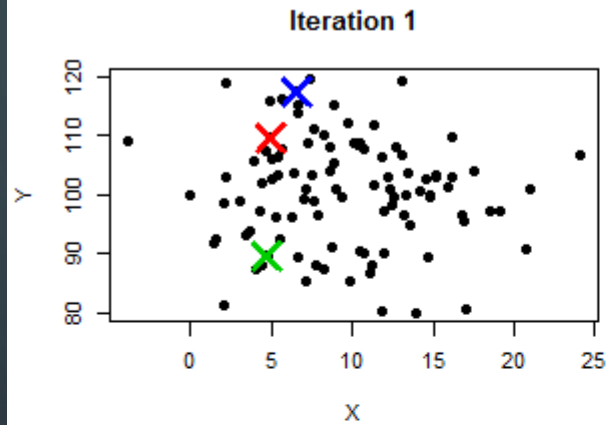
repeat

assign each item t_i to the clusters which has the closest mean;

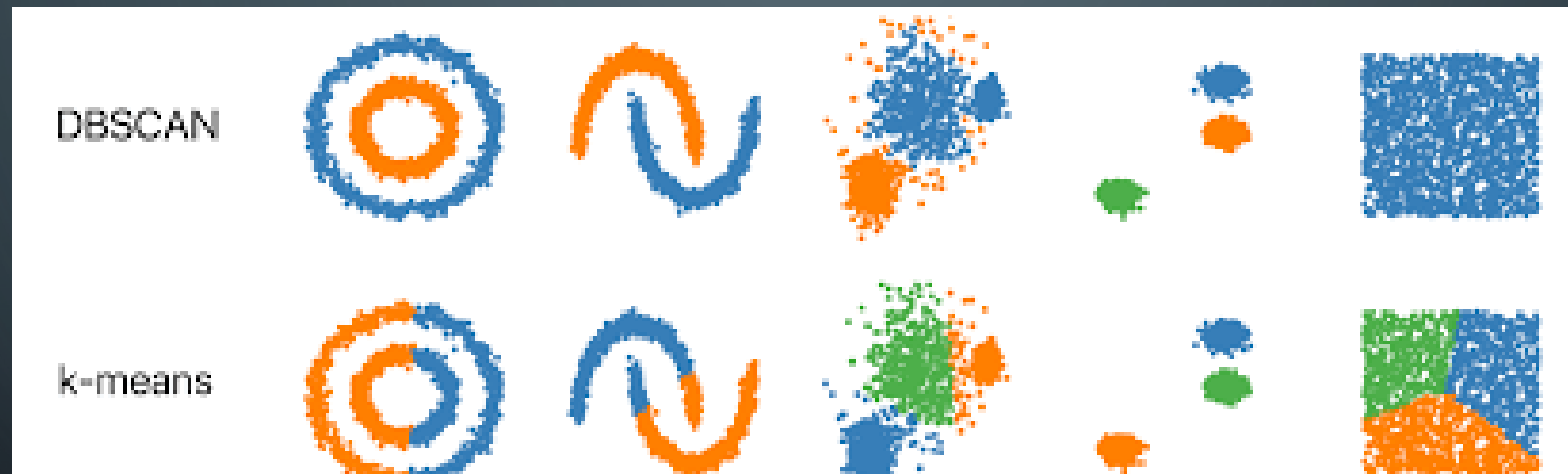
calculate new mean for each cluster;

until convergence criteria is met;

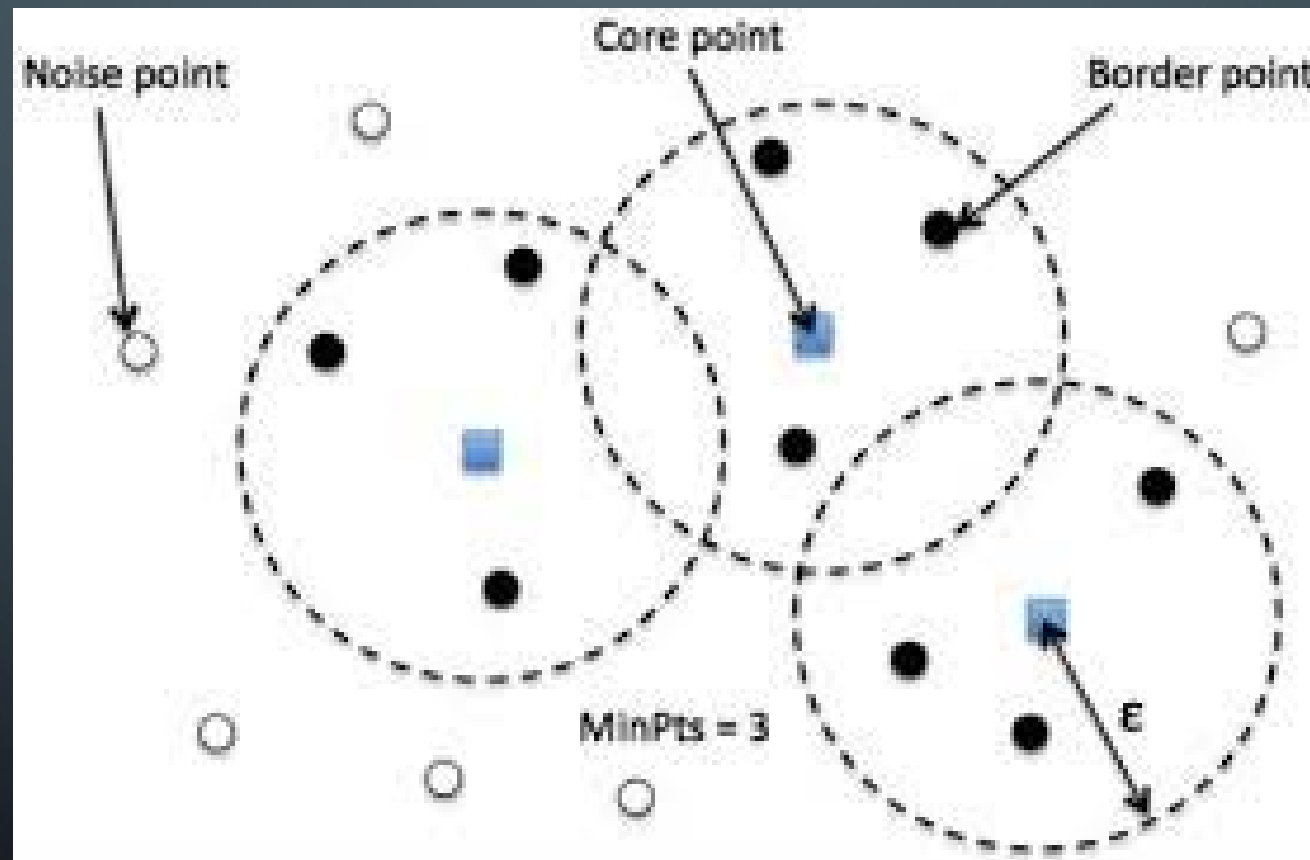
K-MEANS



DB-SCAN



CORE, BORDER AND NOISE POINTS

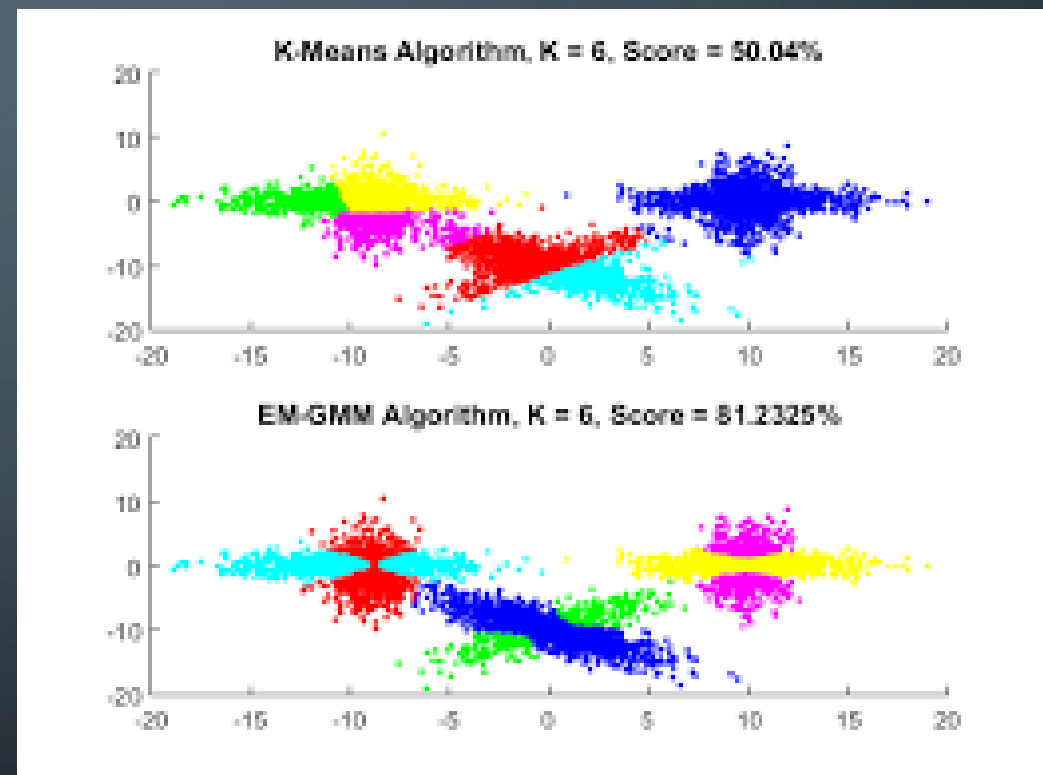
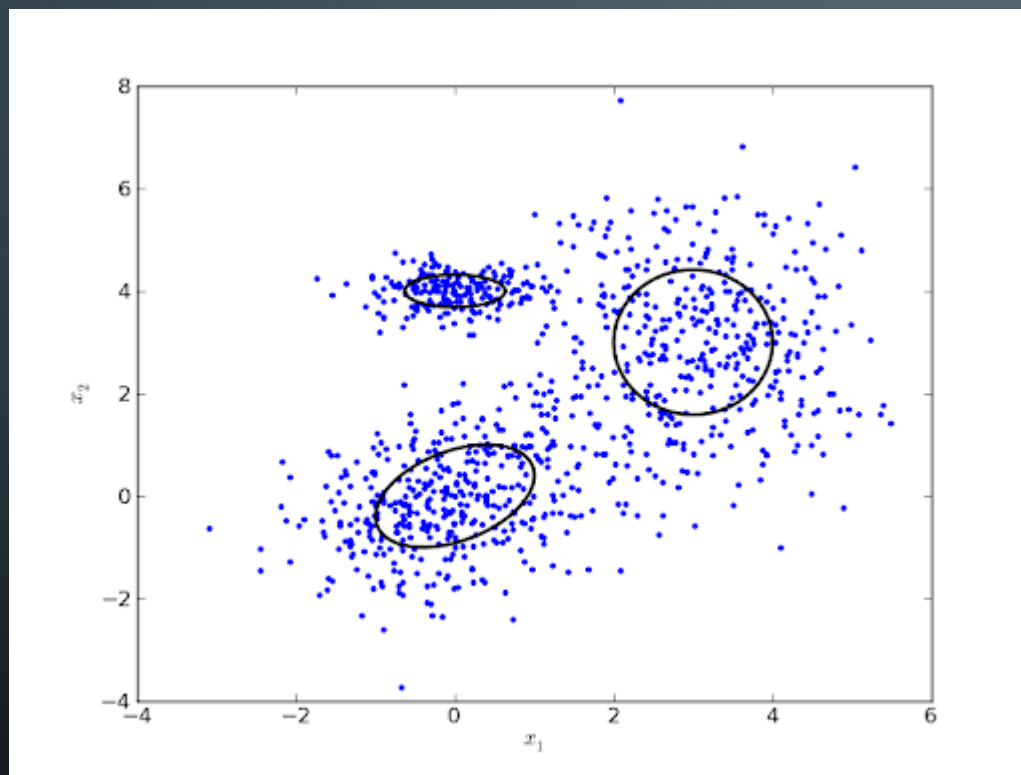


DB-SCAN

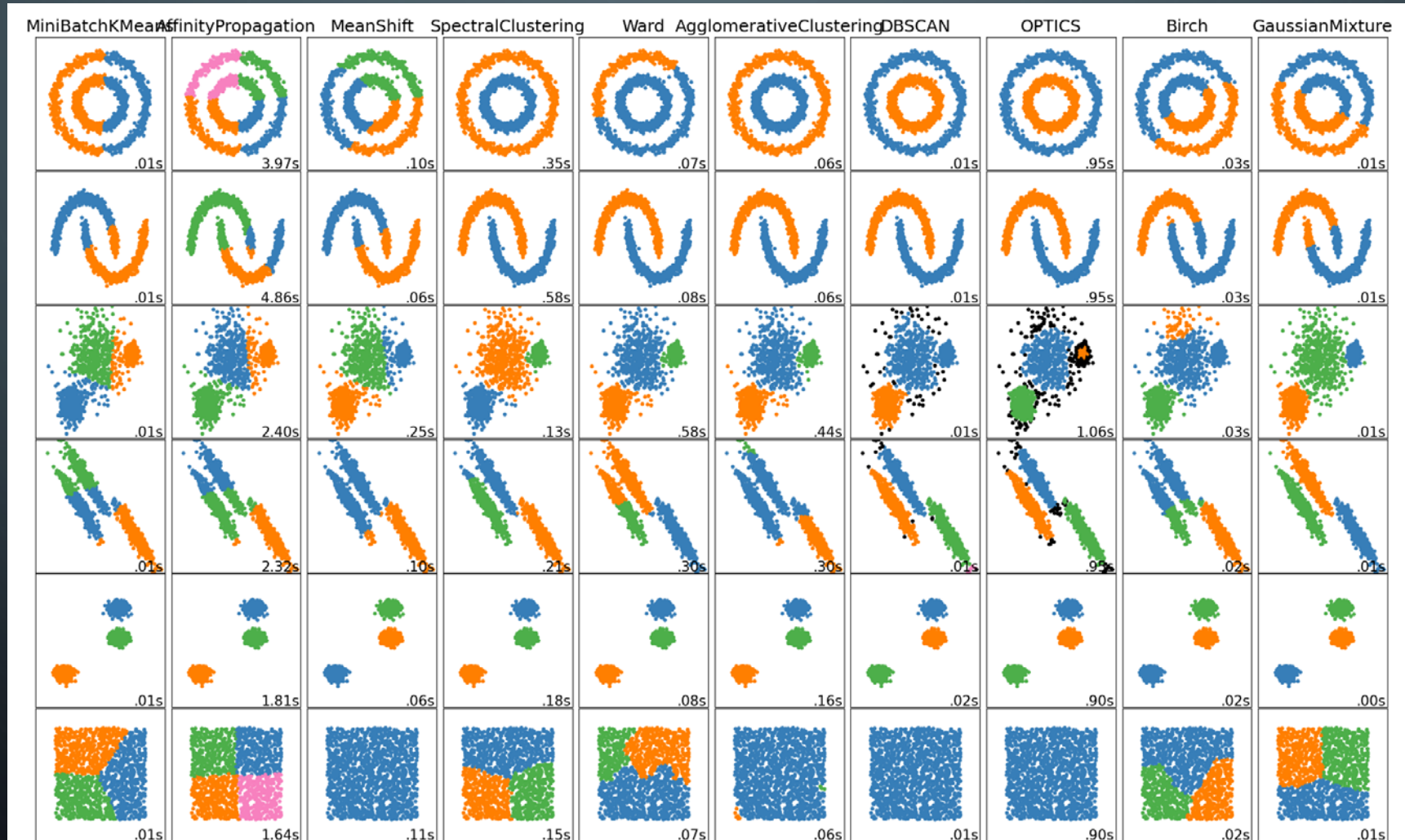
Algorithm 1. DBSCAN algorithm

```
Data:  
Dataset -  $D$ ,  
distance -  $\varepsilon$ ,  
minimum number of points to create dense region -  $minPts$   
1 begin  
2    $C \leftarrow 0$   
3   for each point  $P$  in dataset  $D$  do  
4     if  $P$  is visited then  
5       Continue to next  $P$   
6     end  
7     else  
8       mark  $P$  as visited  
9        $nbrPts \leftarrow$  points in  $\varepsilon$ -neighborhood of  $P$  (distance function)  
10      if  $sizeof(nbrPts) < minPts$  then  
11        mark  $P$  as NOISE  
12      end  
13      else  
14         $C \leftarrow NewCluster$   
15        Call Expand Cluster Function( $P, nbrPts, C, minPts$ )  
16      end  
17    end  
18  end  
19 end
```


GAUSSIAN MIXTURE MODEL



CLUSTERING IN DIFFERENT DATA



The Jupyter logo is centered in the image. It consists of two orange, curved, crescent-like shapes that form a circle around the word "jupyter". The word "jupyter" is written in a white, lowercase, sans-serif font. There are four white circles of varying sizes positioned around the logo: one at the top left, one at the top right, one at the bottom left, and one at the bottom right. The background is a dark blue gradient. In the corners, there are faint, light blue circuit-like patterns with lines and small circles.

jupyter



DIMENSIONALITY REDUCTION

MOHAMMAD GHODDOSI

DIMENSIONALITY REDUCTION

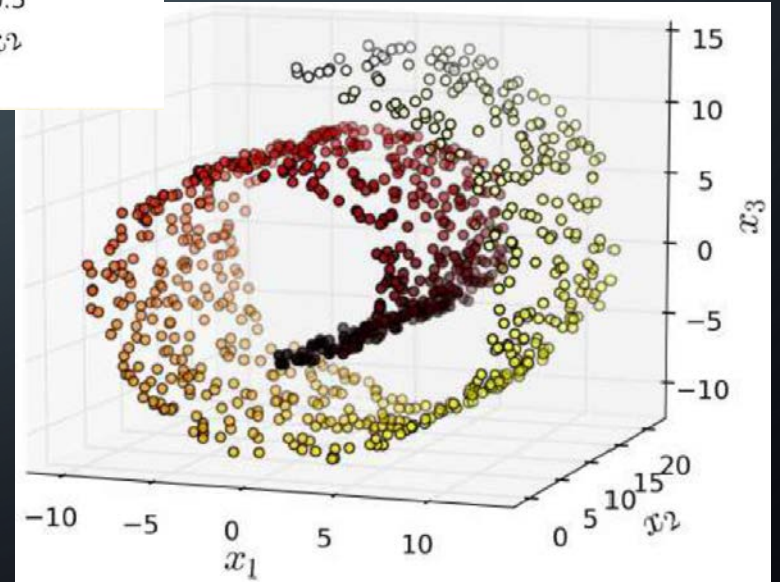
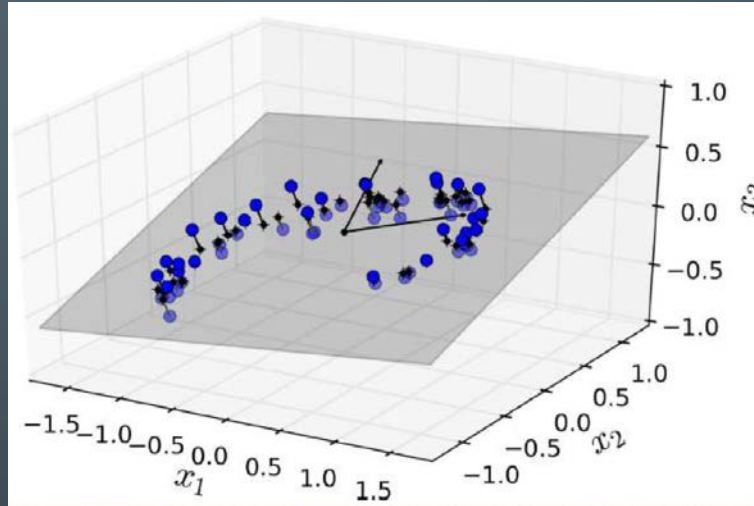
- transformation of data from a high-dimensional space into a low-dimensional space
- low-dimensional representation retains some meaningful properties of the original data

DIMENSIONALITY REDUCTION USE CASES

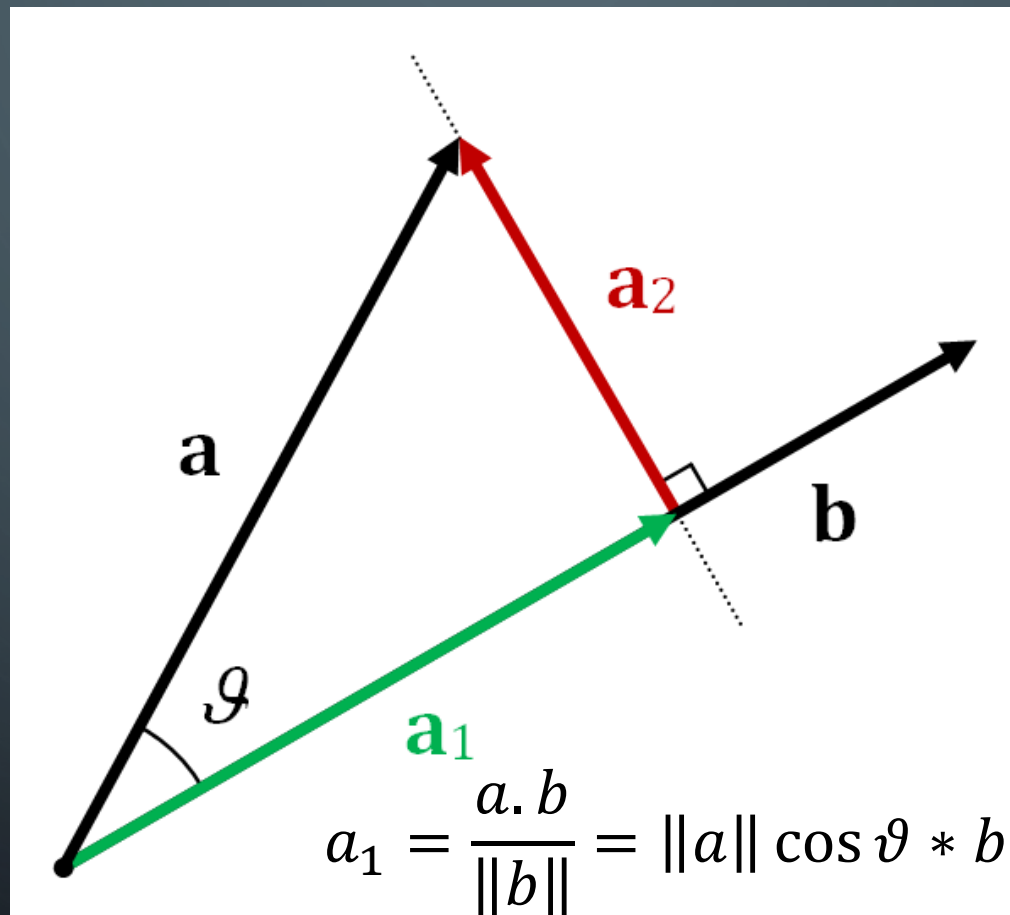
- Prevent overfitting
- Visualization
- Faster and more efficient models
- Storing datasets
- Decorrelation

DIMENSIONALITY REDUCTION TECHNIQUES

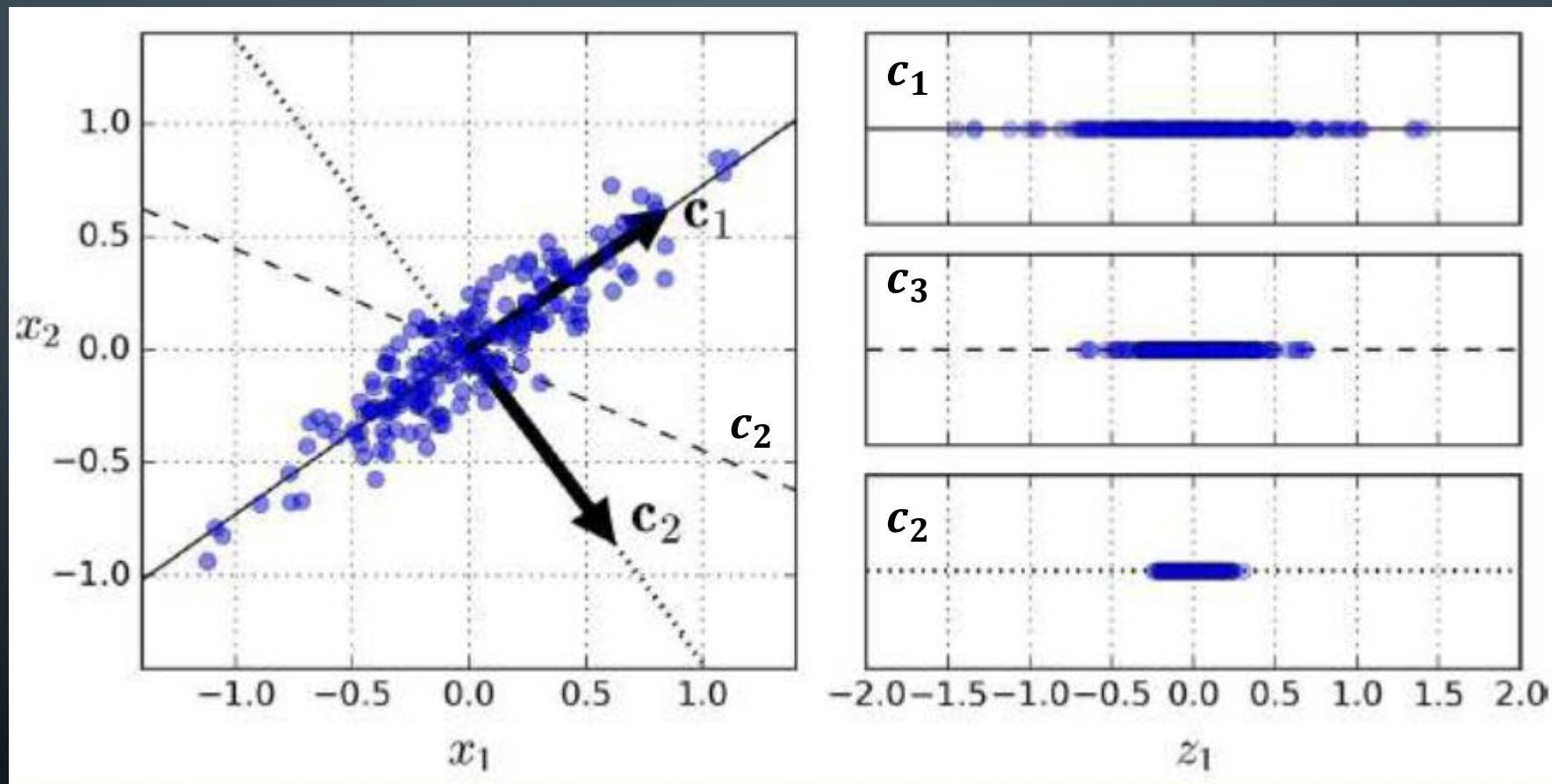
- Feature selection
- Feature projection
 - PCA
- Manifold learning
 - LLE



PROJECTION



PRINCIPLE COMPONENT ANALYSIS (PCA)



PCA IN VECTORS

$$M = U.\Sigma.V^T \text{ (SVD decomposition)}$$

$$U^*U = UU^* = I$$

Σ is diagonal

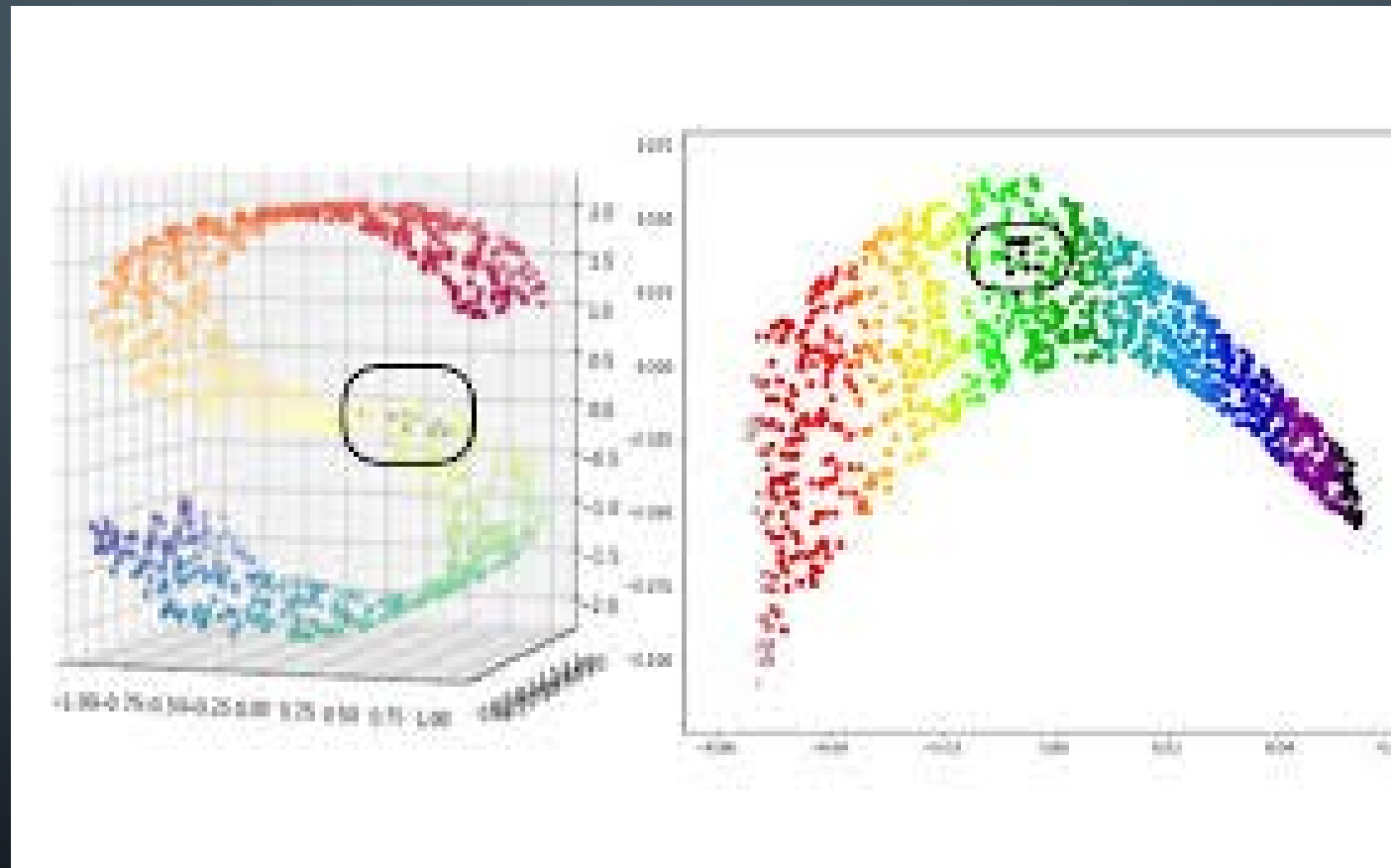
V is components

if M is real, U and V are real and orthogonal

LOCAL LINEAR EMBEDDING (LLE)

- Measuring how each training instance linearly relates to its closest neighbors
- Then looking for low-dimensional representation where these local relationships are best preserved

LLE



LLE ALGORITHM

- Step I: find k-nearest neighbors
- Step II: find W such that:

$$W = \underset{W}{\operatorname{argmin}} \sum_{i=1}^m \left\| X^{(i)} - \sum_{j=1}^m w_{i,j} X^{(j)} \right\|^2$$
$$W = \begin{cases} w_{i,j} & \text{if } X^{(j)} \text{ is not in KNN } X^{(i)} \\ \sum_{j=1}^m w_{i,j} = 1 & \text{for each } i \end{cases}$$

LLE ALGORITHM

- Step III: find Z matrix in lower dimension such that:

$$Z = \operatorname{argmin}_Z \sum_{i=1}^m \left\| Z^{(i)} - \sum_{j=1}^m w_{i,j} Z^{(j)} \right\|^2$$

OTHER TECHNIQUES

- Multidimensional scaling (MDS)
 - Preserve the distances between the instances
- Isomap
 - Forms nearest neighbor graph and try to preserve geodesic distances
- T-SNE
 - Keep similar instances close and dissimilar apart (cluster visualization)
- LDA
 - Classification algorithm, keep classes as far apart as possible.

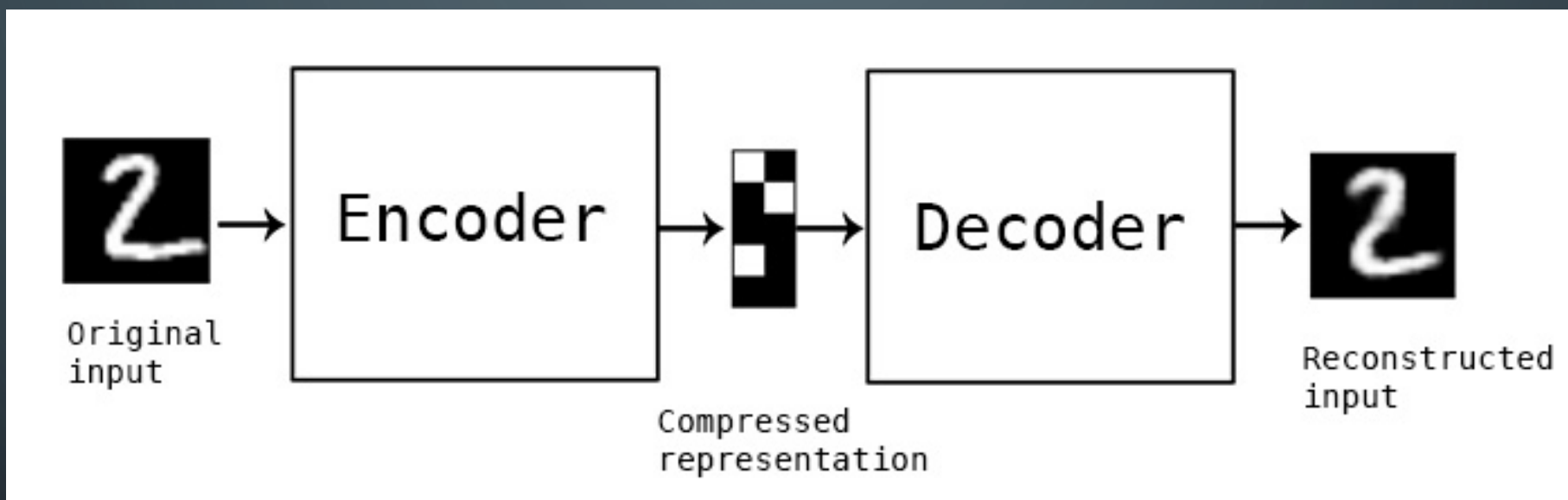
The Jupyter logo is centered in the image. It consists of two orange, curved, crescent-like shapes that form a circle around the word "jupyter". There are four white circles of varying sizes positioned around the logo: one at the top-left, one at the top-right, one at the bottom-left, and one at the bottom-right. The word "jupyter" is written in a white, lowercase, sans-serif font.

jupyter

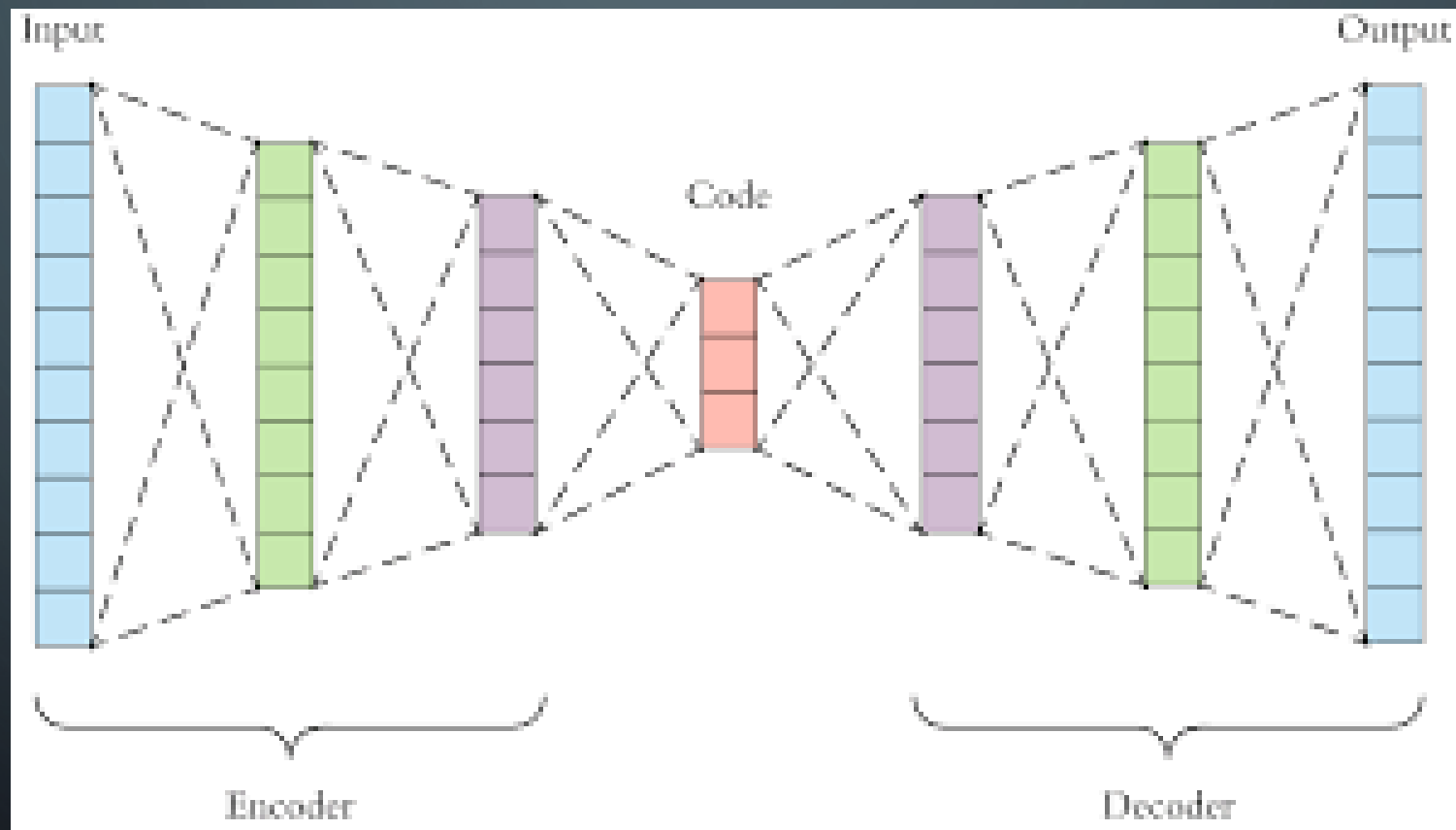
AUTOENCODERS

- Artificial Neural Network
- learns efficient data representation in an unsupervised manner
- Used in dimensional reduction

AUTOENCODERS, IDEA



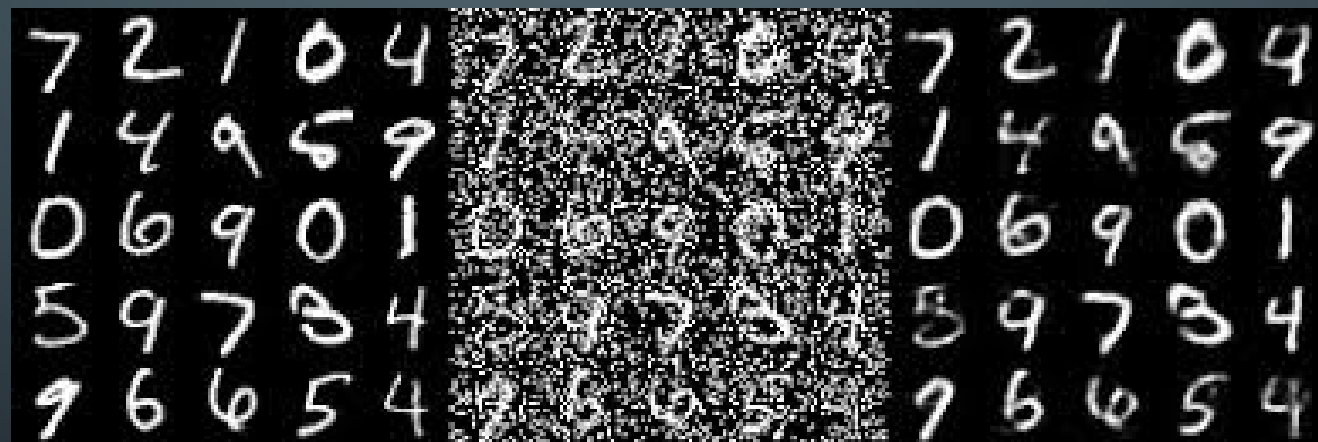
UNDERCOMPLETE AUTOENCODERS



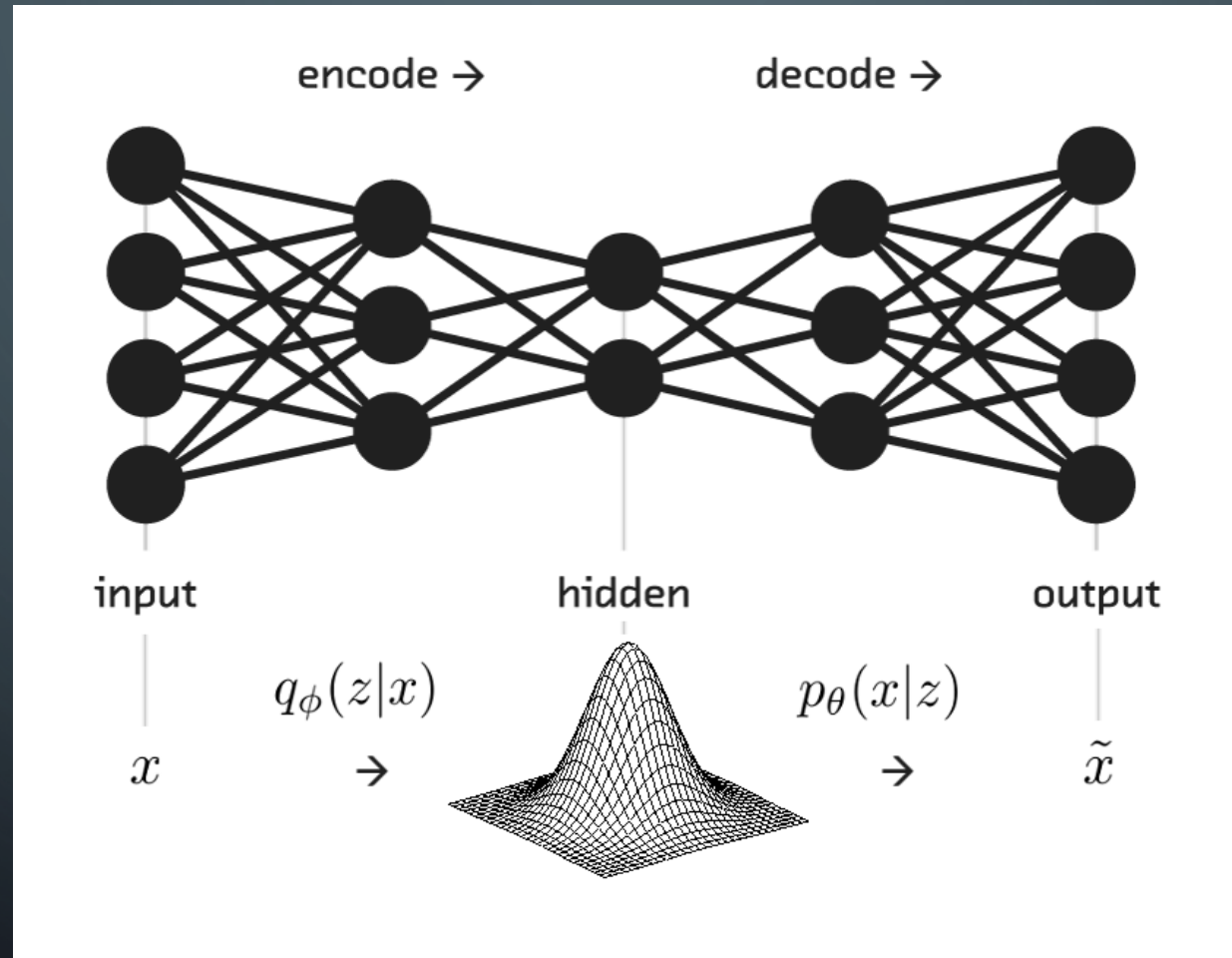
AUTOENCODERS VARIATIONS

- undercomplete Autoencoders
 - bottleneck
- Sparse Autoencoders
 - L1 regularization
 - Better for classification
- Denoising Autoencoders
 - Input is corrupted
- Variational Autoencoders
 - generates new samples

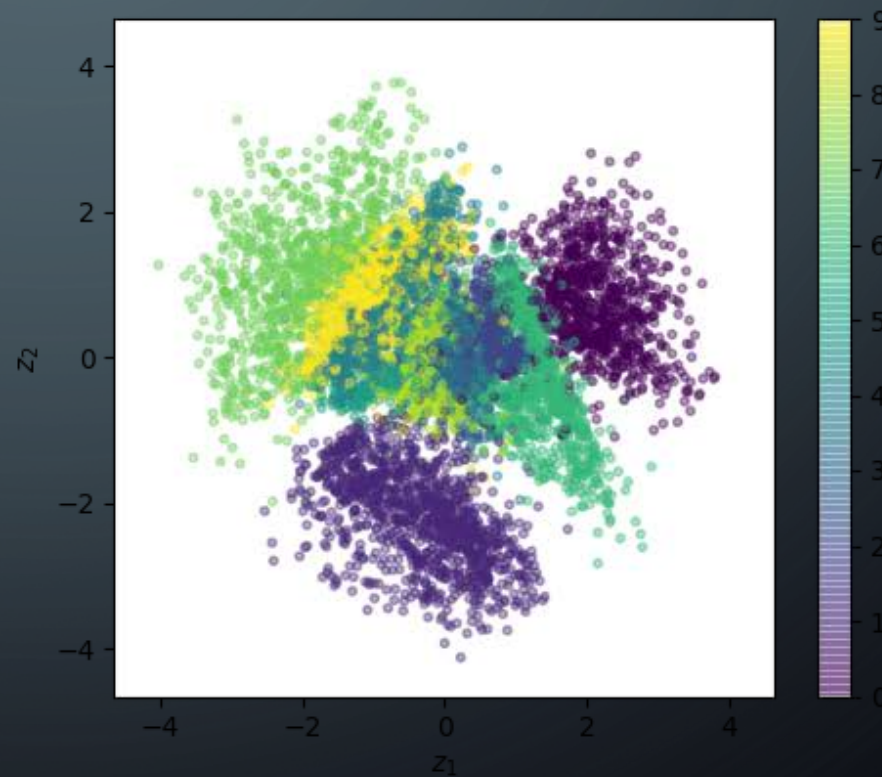
DENOISING AUTOENCODERS



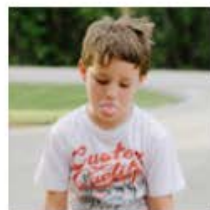
VARIATIONAL AUTOENCODERS



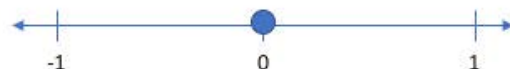
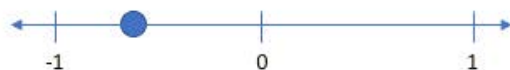
VARIATIONAL AUTOENCODERS (MNIST)



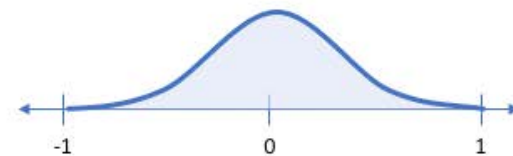
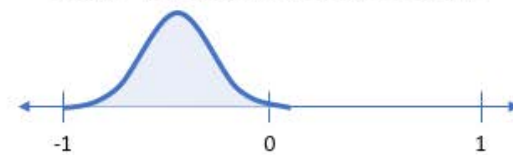
VA IDEA



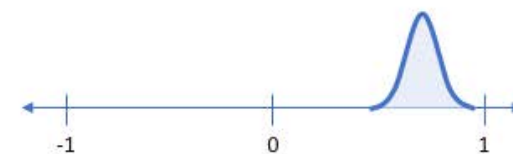
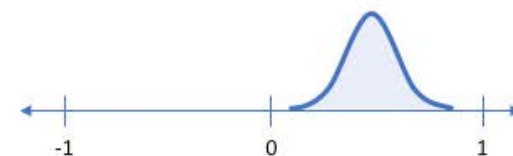
Smile (discrete value)



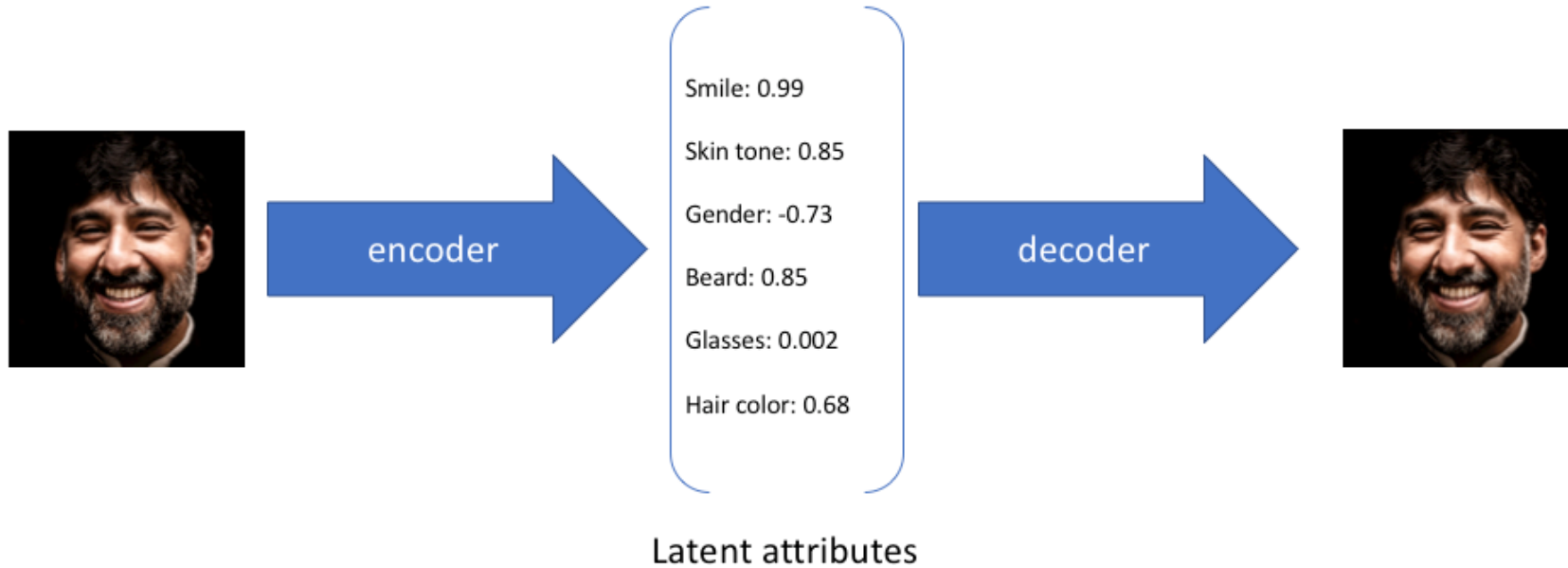
Smile (probability distribution)



vs.



NORMAL AUTOENCODERS



VARIATIONAL AUTOENCODERS

