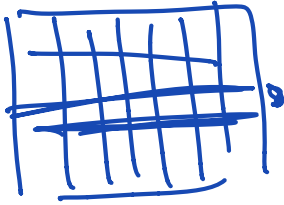


Unsupervised Learning (2)

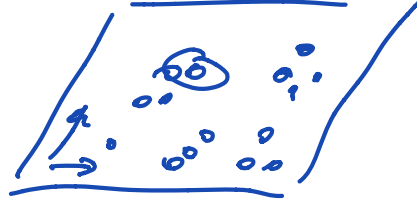
Machine Learning for Data Science 1

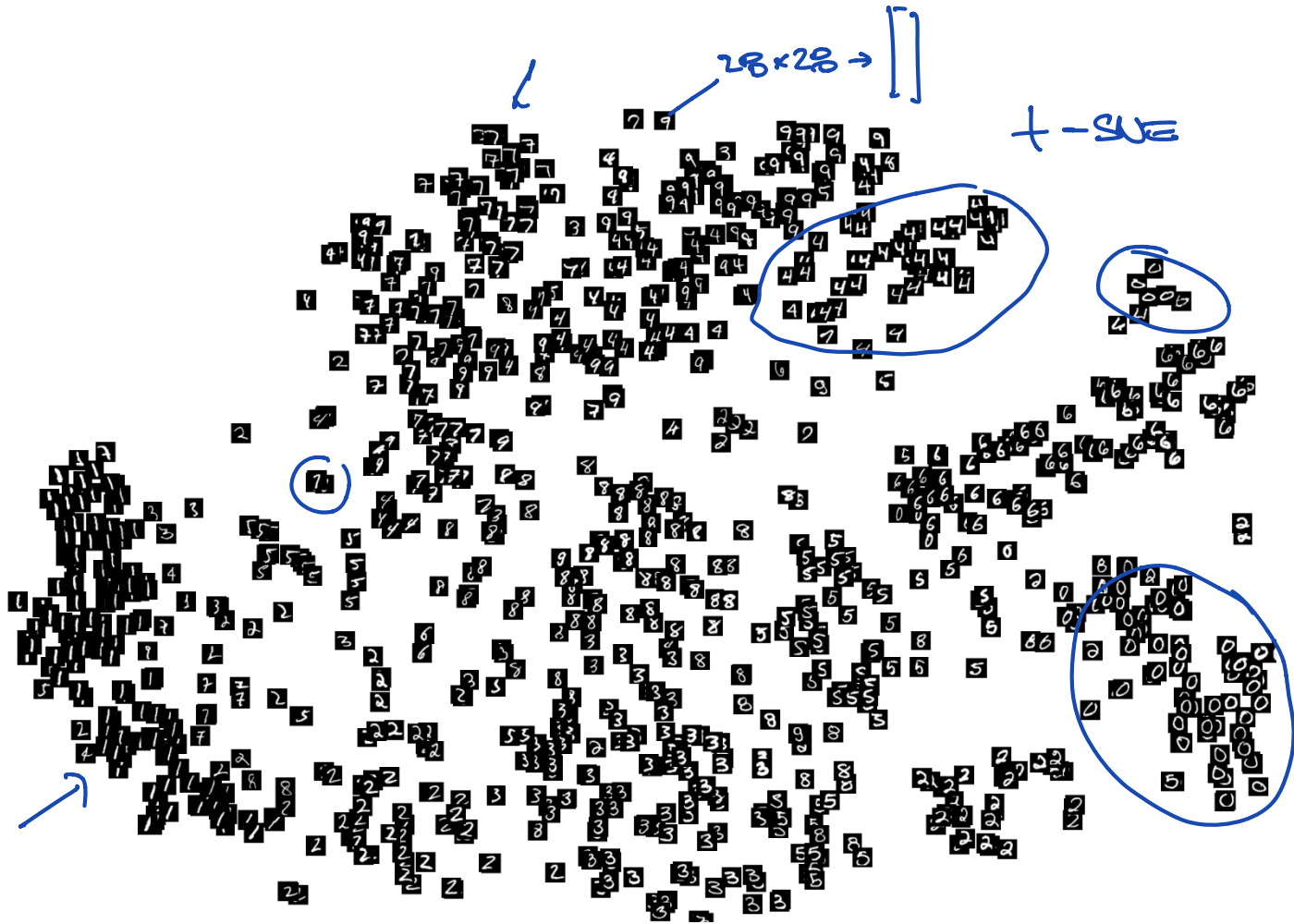
Latent spaces :

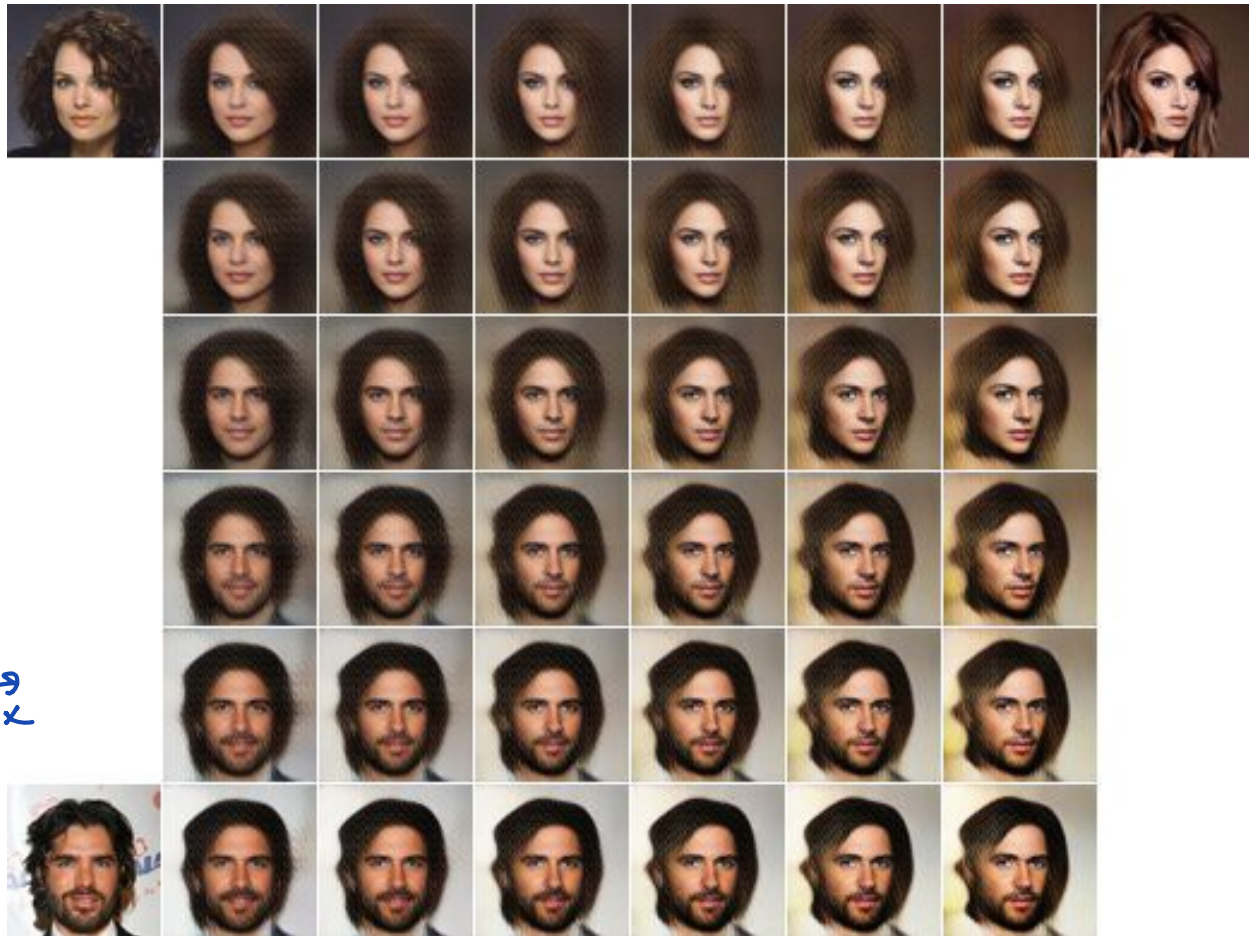
representation of compressed data
in which similar data
points (from our input) \rightarrow go & see
reside together



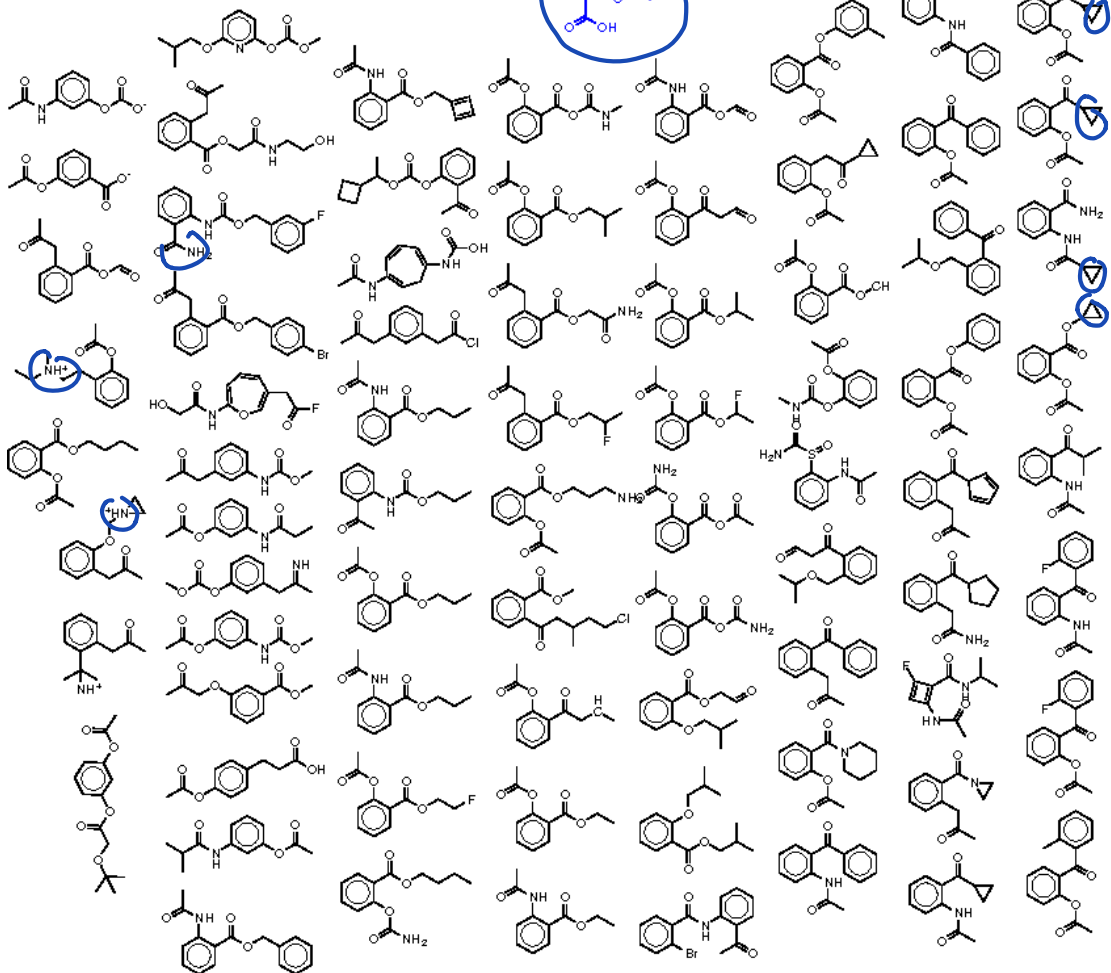
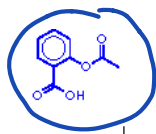
10, 100
1000's ..







lokut syare ej chemicals



PCA

MDS

t-SNE

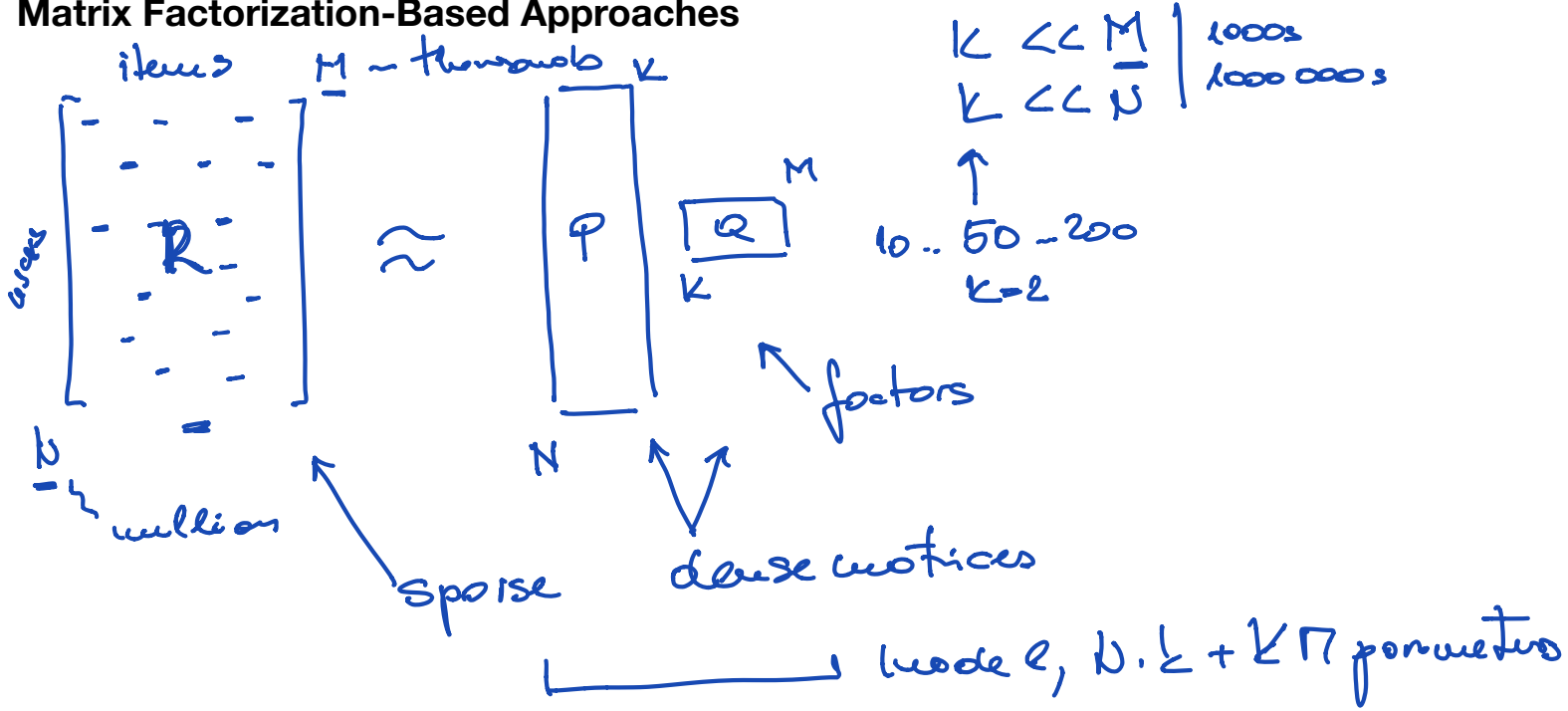
dimensionality reduction
construct latent spaces

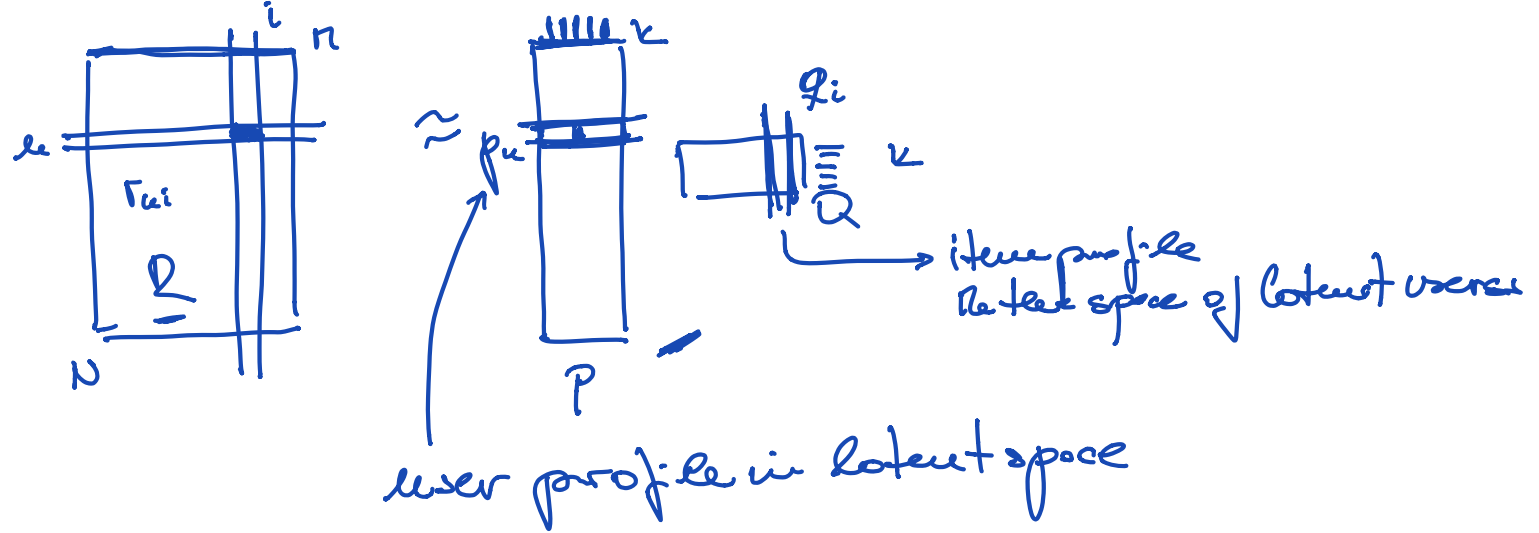
PCA - projection

MDS, t-SNE - embedding

Sparse Data,

Matrix Factorization-Based Approaches





$$r_{ui} \approx p_u^T q_i$$

$$\hat{r}_{ui} \approx r_{ui}$$

$$\hat{r}_{ui} = p_u^T q_i = \sum_{k=1}^k p_{uk} \cdot q_{ki}$$

$$e_{ui}^2 = (r_{ui} - \hat{r}_{ui})^2$$

$$= \left(r_{ui} - \sum_k p_{uk} \cdot q_{ki} \right)^2$$

user user

gradient (stochastic gradient descent)

$$\frac{\partial \mathcal{L}_{ui}^2}{\partial p_{uk}} = -2(r_{ui} - \hat{r}_{ui}) g_{ki} \quad \text{for } k=1 \dots K$$
$$= -2 e_{ui} \cdot g_{ki}$$

$$\frac{\partial \mathcal{L}_{ui}^2}{\partial g_{ki}} = -2(r_{ui} - \hat{r}_{ui}) p_{uk}$$
$$= -2 e_{ui} p_{uk}$$

update rules

$$p_{uk} \leftarrow p_{uk} - \alpha \frac{\partial \mathcal{L}_{ui}^2}{\partial p_{uk}} = p_{uk} + \alpha e_{ui} g_{ki} \quad \text{for } k=1 \dots K$$
$$g_{ki} \leftarrow g_{ki} + \alpha e_{ui} p_{uk} \quad \text{for } k=1 \dots K$$

regularization $\|P\|_F, \|Q\|_F$ $J = e_{ui}^2 + \lambda \sum p_u^T p_u + \lambda \sum q_i^T q_i$
initialization $P, Q: [0, 0.001]$ stopping criteria RMSE

Tekoc et al
 IJIR 2008

Input: \mathcal{T}' : training set, η : learning rate, λ : regularization factor
Output: $\mathbf{P}^*, \mathbf{Q}^*$: the user and item feature matrices

- 1 Partition \mathcal{T}' into two sets: $\mathcal{T}'_I, \mathcal{T}'_{II}$ (validation set)
- 2 Initialize \mathbf{P} and \mathbf{Q} with small random numbers.
- 3 **loop** until the terminal condition is met. One epoch:
 - 4 iterate over each (u, i, r_{ui}) element of \mathcal{T}'_I : r_{ui}
 - 5 compute e'_{ui} ;
 - 6 compute the gradient of e'_{ui} , according to Eq. (6);
 - 7 for each k
 - 8 update \mathbf{p}_u , the u -th row of \mathbf{P} ,
 - 9 and \mathbf{q}_i , the i -th column of \mathbf{Q} according to Eq. (7);
 - 10 calculate the RMSE on \mathcal{T}'_{II} ;
 - 11 if the RMSE on \mathcal{T}'_{II} was better than in any previous epoch:
 - 12 Let $\mathbf{P}^* = \mathbf{P}$ and $\mathbf{Q}^* = \mathbf{Q}$.
 - 13 terminal condition: RMSE on \mathcal{T}'_{II} does not decrease during two epochs.
- 14 **end**

Algorithm 1: Training algorithm for RISMF

Oct 2006

circumstances

4200 000 users

18 000 movies

100 420 000 ratings

10%, best customer

June 2009



March 12, 2010

concerned by Federal Trade Commission
privacy concerns

Non-Negative Matrix Factorization

letters to nature

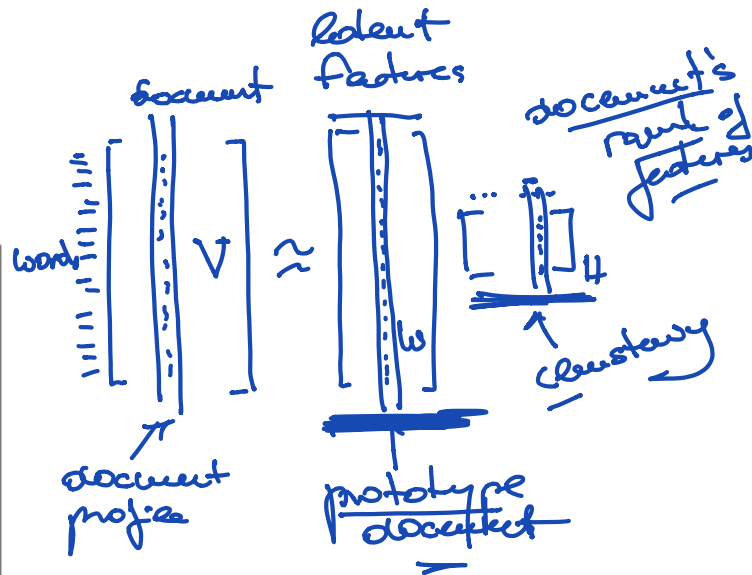
Learning the parts of objects by non-negative matrix factorization

Daniel D. Lee* & H. Sebastian Seung*†

* Bell Laboratories, Lucent Technologies, Murray Hill, New Jersey 07974, USA

† Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

Is perception of the whole based on perception of its parts? There is psychological¹ and physiological^{2,3} evidence for parts-based representations in the brain, and certain computational theories of object recognition rely on such representations^{4,5}. But little is known about how brains or computers might learn the parts of objects. Here we demonstrate an algorithm for non-negative matrix factorization that is able to learn parts of faces and semantic features of text. This is in contrast to other methods, such as principal components analysis and vector quantization, that learn holistic, not parts-based, representations. Non-negative matrix factorization is distinguished from the other methods by its use of non-negativity constraints. These constraints lead to a parts-based representation because they allow only additive, not subtractive, combinations. When non-negative matrix factorization is implemented as a neural network, parts-based representations emerge by virtue of two properties: the firing rates of neurons are never negative and synaptic strengths do not change sign.



$$\min \|V - WH\|_F$$

subj. to $W \geq 0$
 $V \geq 0$

Nature 1999

multiplicative update rules

$$V \sim \frac{WH}{\sum_i} \quad \begin{array}{l} \nearrow \text{weight} \\ \nearrow \text{latent matrices} \end{array}$$

$$H_{ij} \leftarrow H_{ij} \cdot \frac{W^T V_{ij}}{W^T W H_{ij}}$$

$$W_{ij} \leftarrow W_{ij} \cdot \frac{V H^T_{ij}}{W H H^T_{ij}}$$

initialize W, H
repeat
update rules

until convergence
validation
matrix

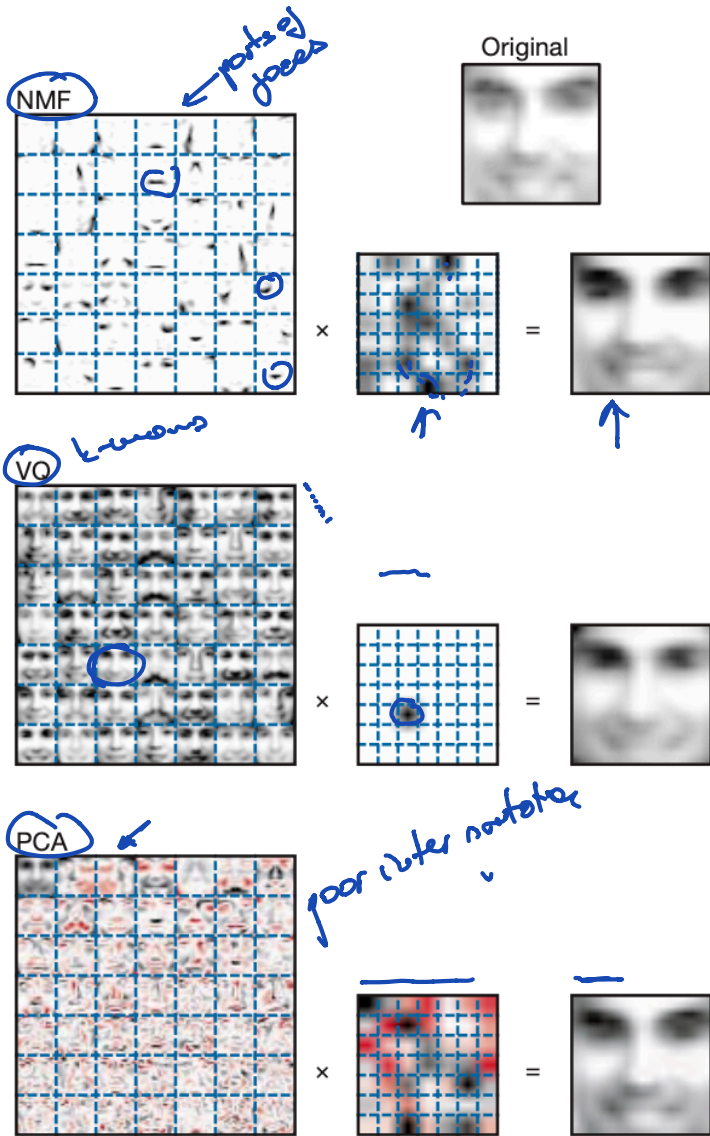
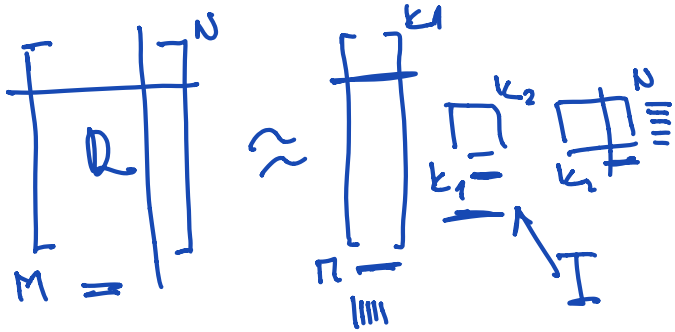


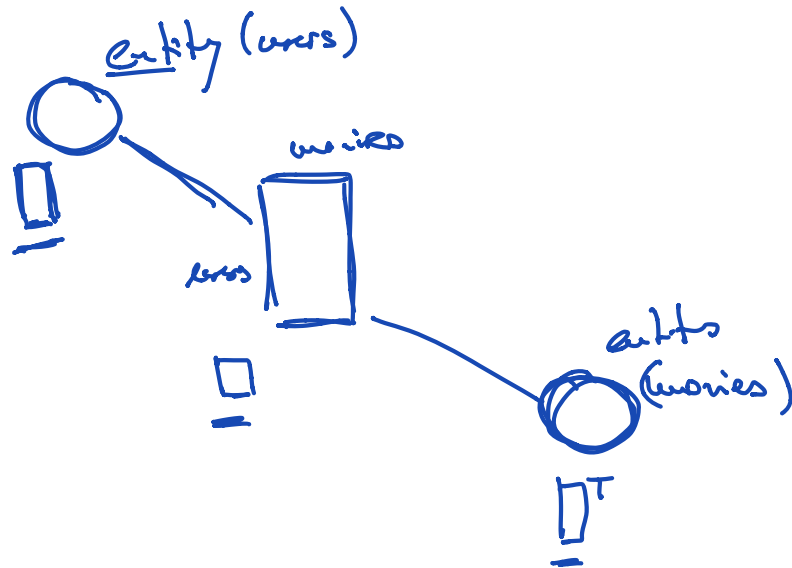
Figure 1 Non-negative matrix factorization (NMF) learns a parts-based representation of faces, whereas vector quantization (VQ) and principal components analysis (PCA) learn holistic representations. The three learning methods were applied to a database of $m = 2,429$ facial images, each consisting of $n = 19 \times 19$ pixels, and constituting an $n \times m$ matrix V . All three find approximate factorizations of the form $V \approx WH$, but with three different types of constraints on W and H , as described more fully in the main text and methods. As shown in the 7×7 montages, each method has learned a set of $r = 49$ basis images. Positive values are illustrated with black pixels and negative values with red pixels. A particular instance of a face, shown at top right, is approximately represented by a linear superposition of basis images. The coefficients of the linear superposition are shown next to each montage, in a 7×7 grid, and the resulting superpositions are shown on the other side of the equality sign. Unlike VQ and PCA, NMF learns to represent faces with a set of basis images resembling parts of faces.

Matrix Tri-Factorization



3 latent factors

$$\underline{k_1} = \underline{k_2}$$





Gene Ontology terms

0	0	1	1	0	0
0	1	0	1	0	0
0	1	1	0	0	0
1	0	0	1	1	0
0	1	1	1	0	1



Approximation of gene-to-function matrix

0.16	0.56	0.55	0.75	0.16	0.26
0.18	0.55	0.54	0.77	0.18	0.26
-0.25	0.76	0.78	0.34	-0.25	0.34
0.86	-0.05	-0.12	1.19	0.86	0.01
0.00	1.10	1.10	1.02	0.00	0.51

Gene recipe matrix

0.57	0.55
0.57	0.58
0.51	0.08
0.42	1.20
0.94	0.64

Backbone matrix

-2.94	3.05
2.05	-1.27

Function recipe matrix

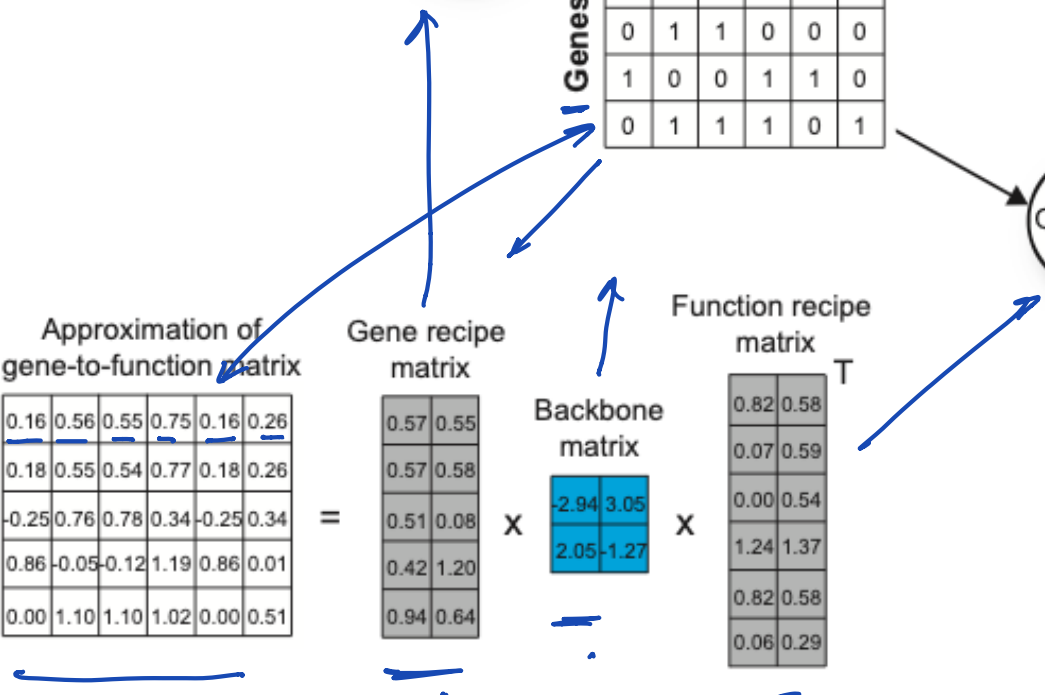
0.82	0.58
0.07	0.59
0.00	0.54
1.24	1.37
0.82	0.58
0.06	0.29

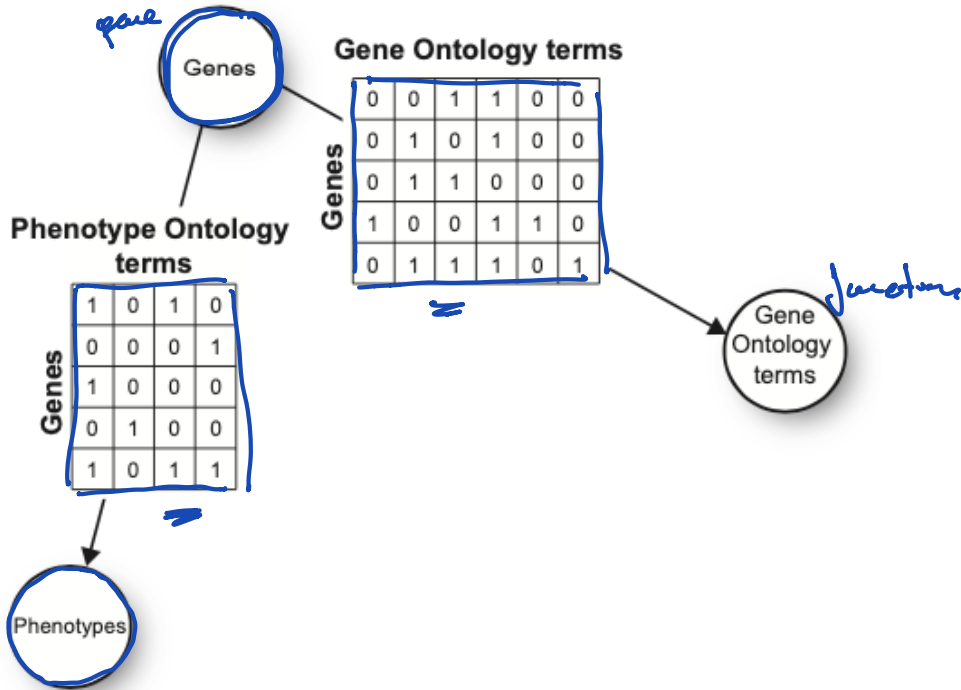
=

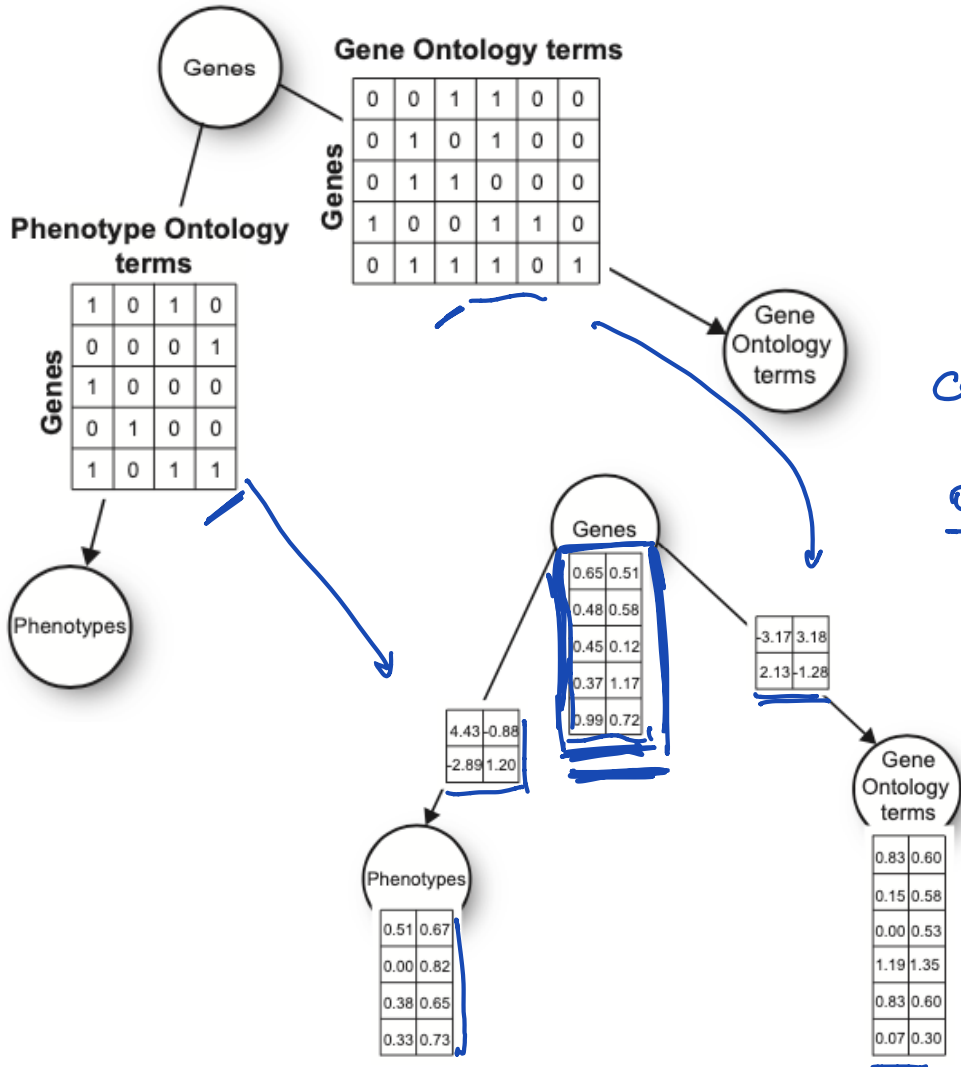
X

X

T

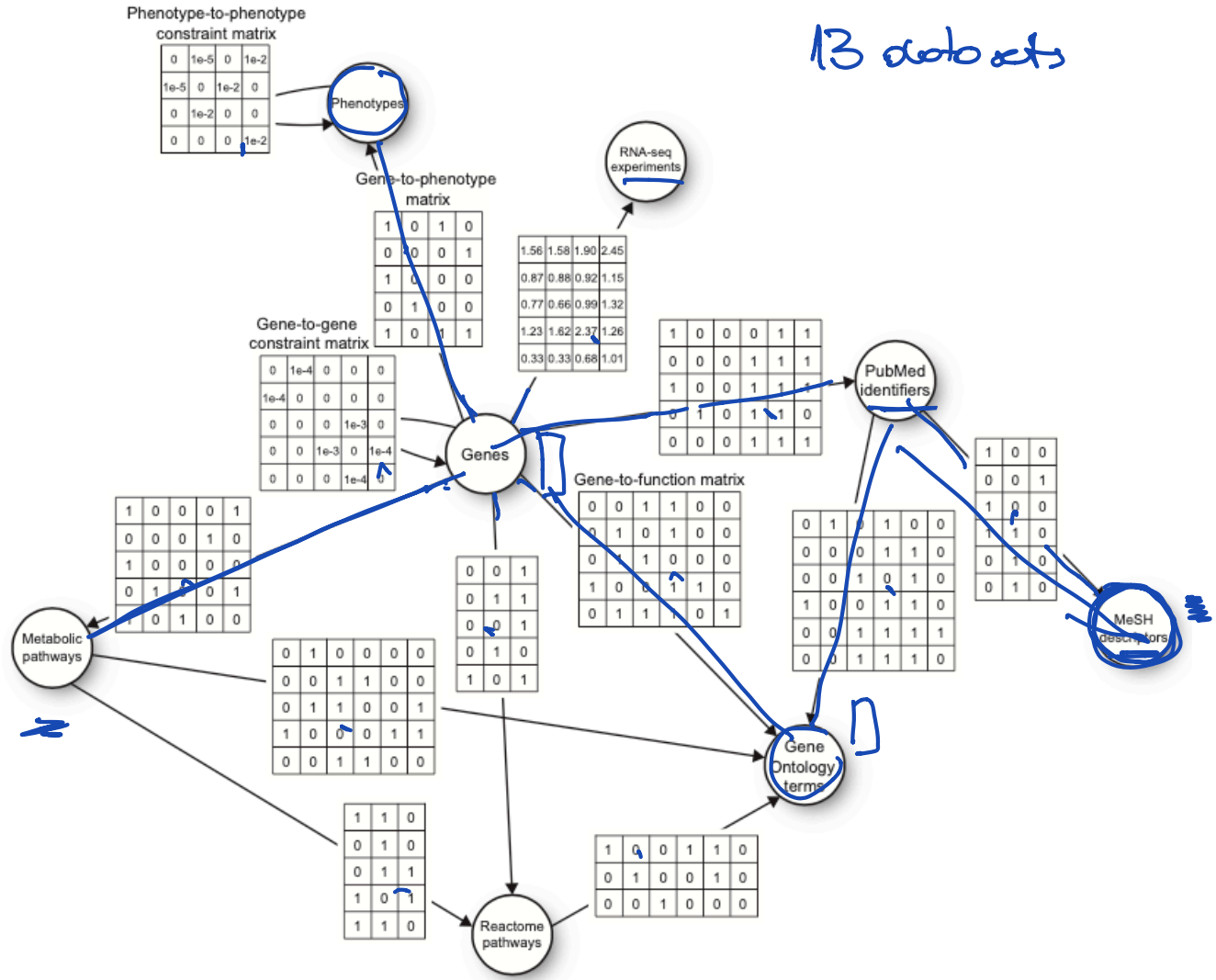




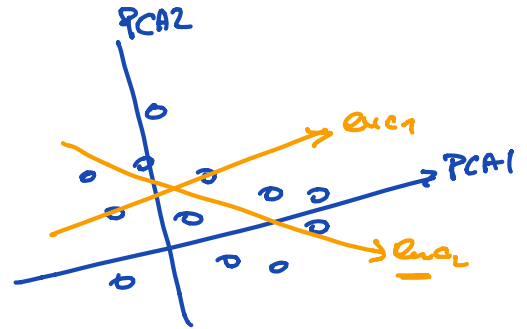
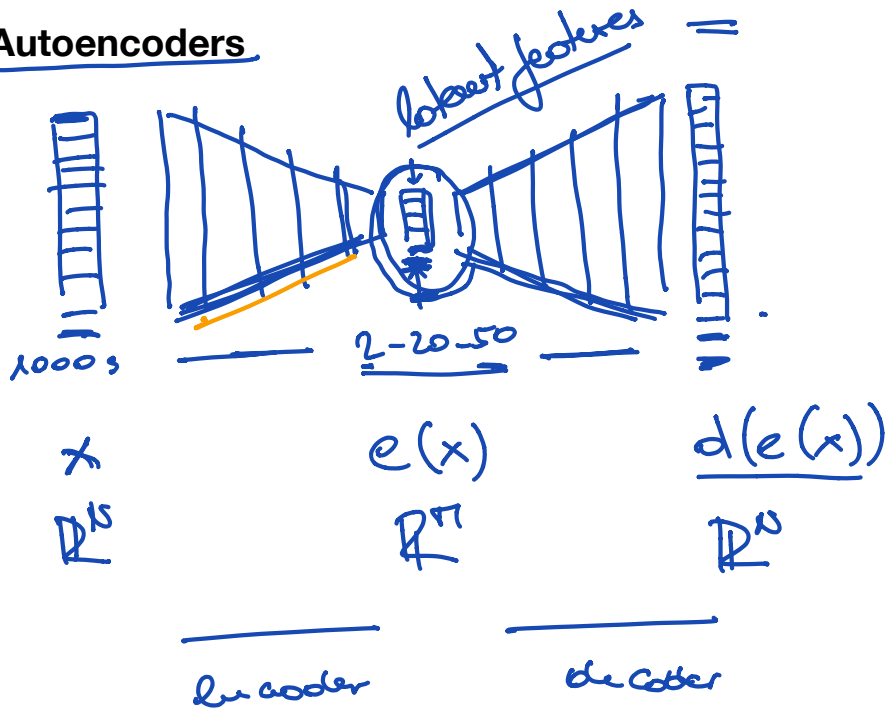


collective best factorization
data fusion

13 datasets



Autoencoders

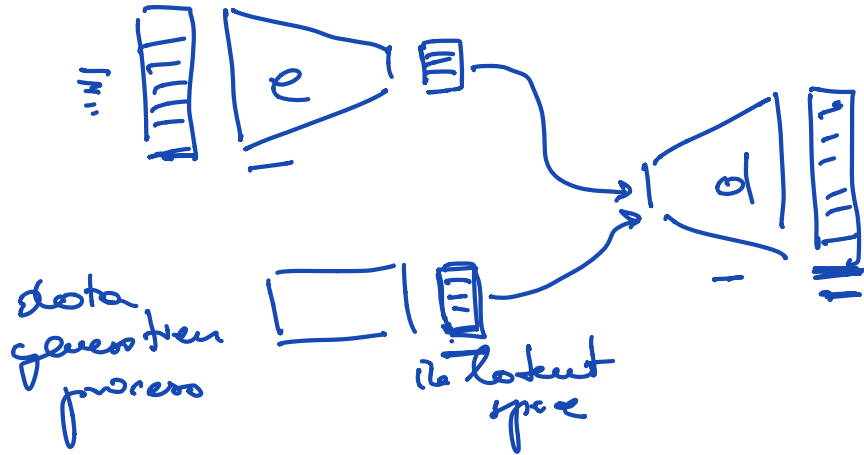


$$n \ll N$$

$$\min_x \|x - d(e(x))\|_F$$

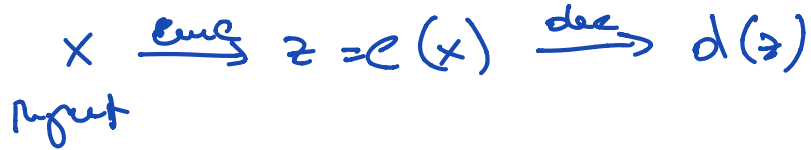
Variational Autoencoders

- no regularity of latent space
- generate new data

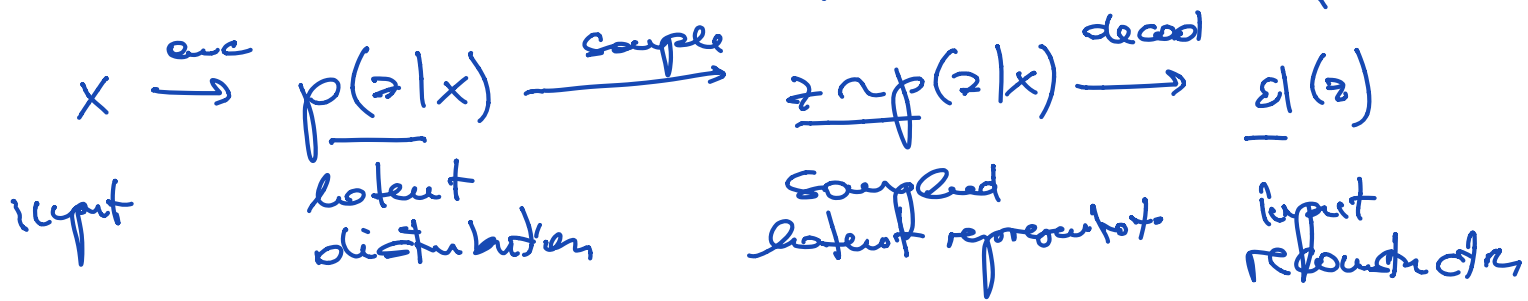


- training is regularized
- latent space has "nice" properties in terms of density
 - ↳ enables generative process

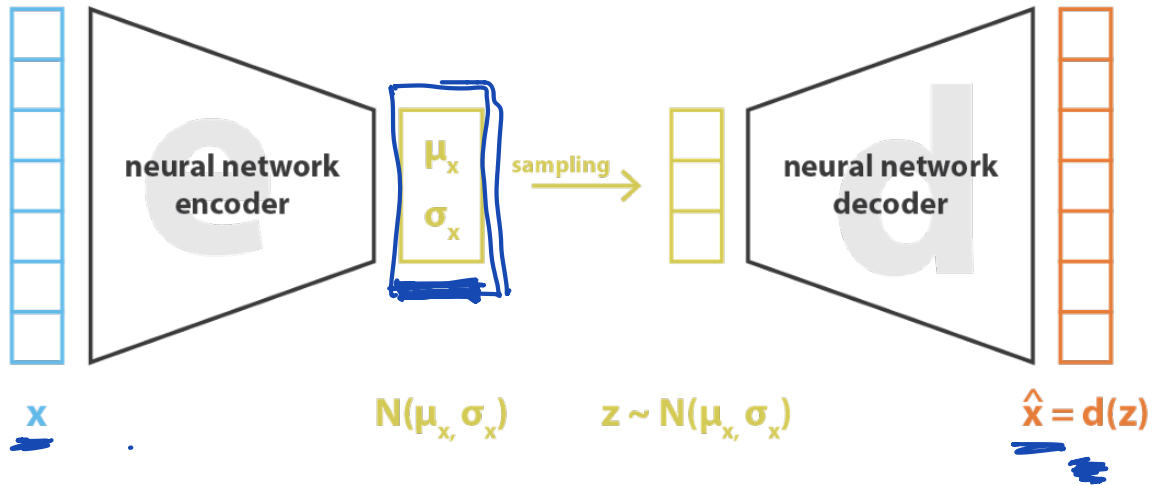
autoencoder (deterministic)



variational autoencoder (probabilistic)



In practice: learned distributions
encoder is trained to
return the mean and cov. matr.



$$\text{loss} = \underbrace{\|x - \hat{x}\|^2}_{\text{reconstruction loss}} + \underbrace{\text{KL}[N(\mu_x, \sigma_x), N(0, I)]}_{\text{KL divergence}} = \|x - d(z)\|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)]$$

Clustering

grouping of objects
objects in the same group
are more similar to
each other than to
objects in other groups

Connectivity-Based Clustering

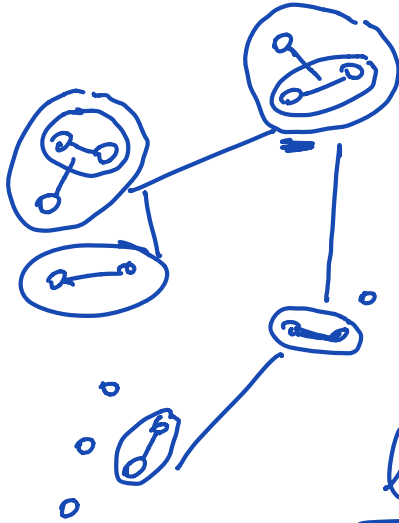
→ hierarchical clustering

similarity

distance

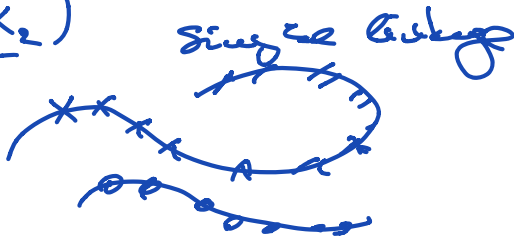
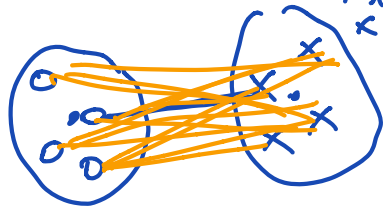
distances:

between the objects



linkage

$$D(c_1, c_2) = \min_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2)$$



$$D(c_1, c_2) = \max_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2)$$

complete linkage



$$D(c_1, c_2) = \frac{1}{|c_1| |c_2|} \sum_{x_1 \in c_1} \sum_{x_2 \in c_2} D(x_1, x_2)$$

average linkage

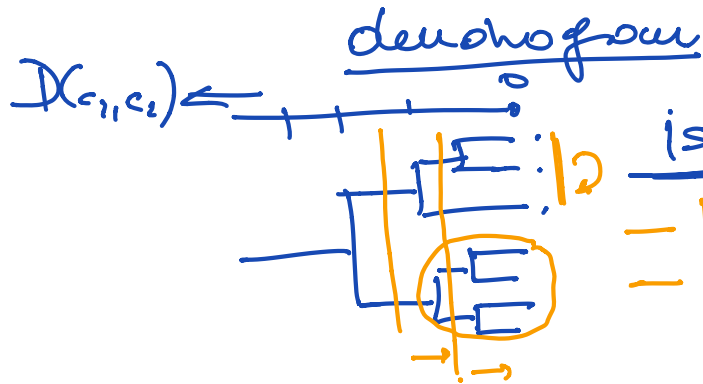
$$\underline{\underline{D(c_1, c_2) = \sum_{x \in c_1 \cup c_2} D(x, \mu_{c_1 \cup c_2})}}$$

Word linkage

↓
minimizes the variance within the clusters

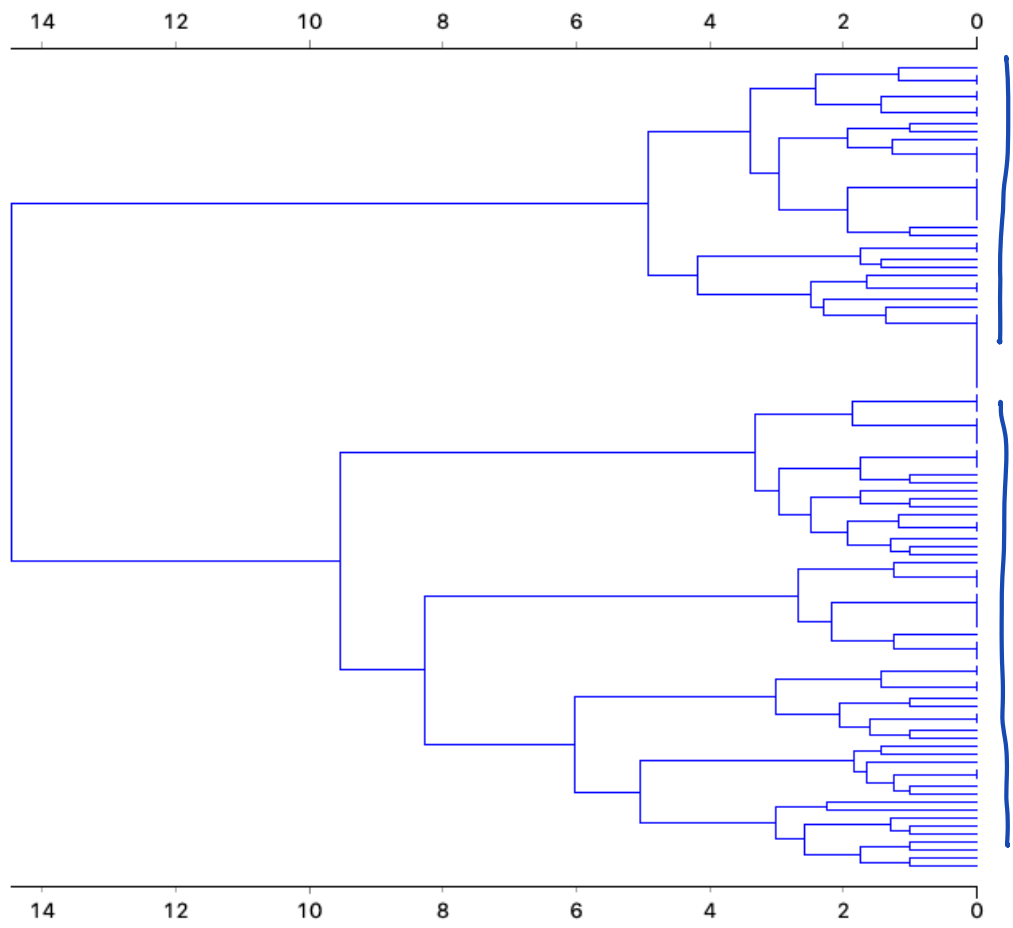
Hierarchical clustering

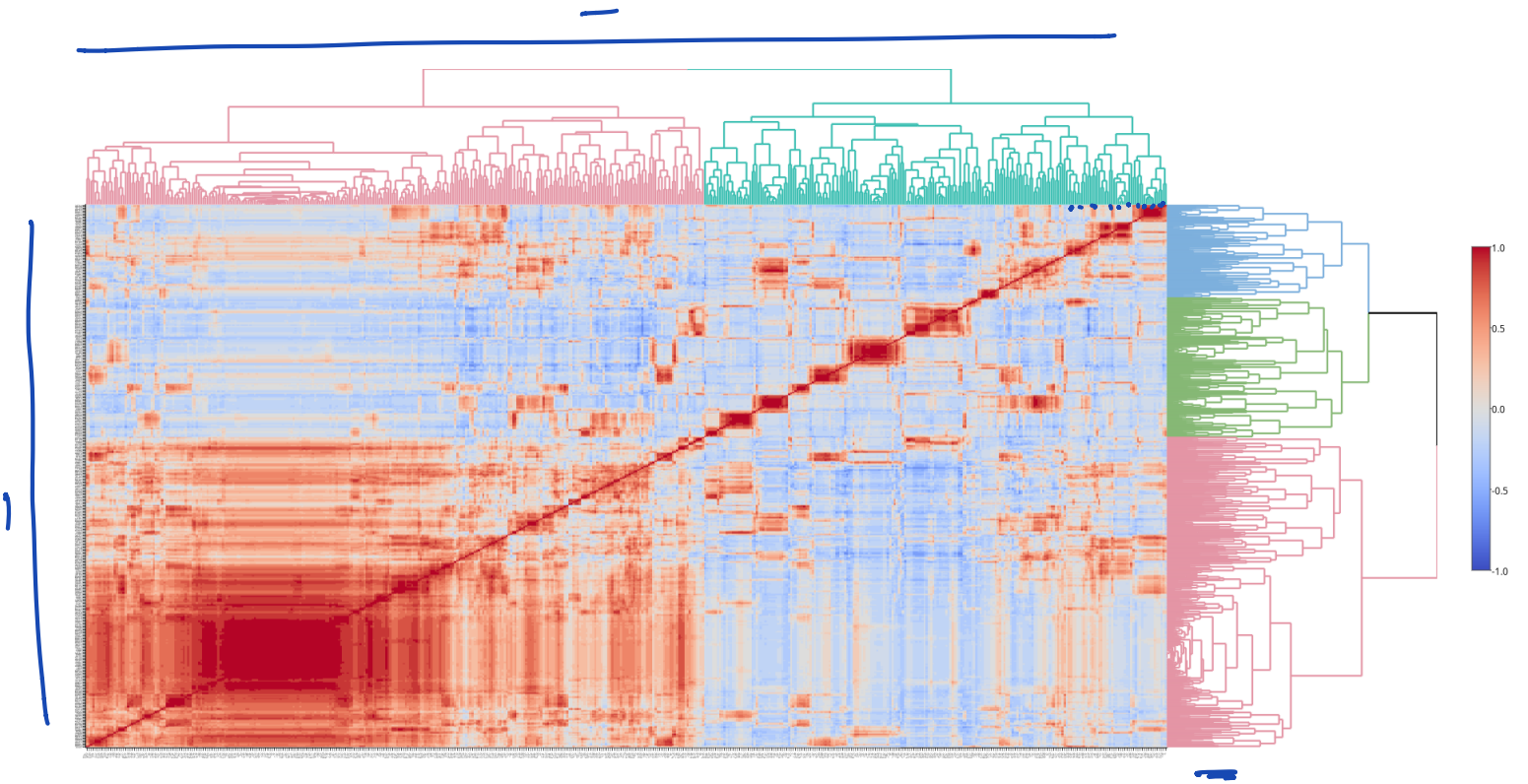
- each object is its own cluster
- merge two closest clusters
- repeat, until just 1 cluster

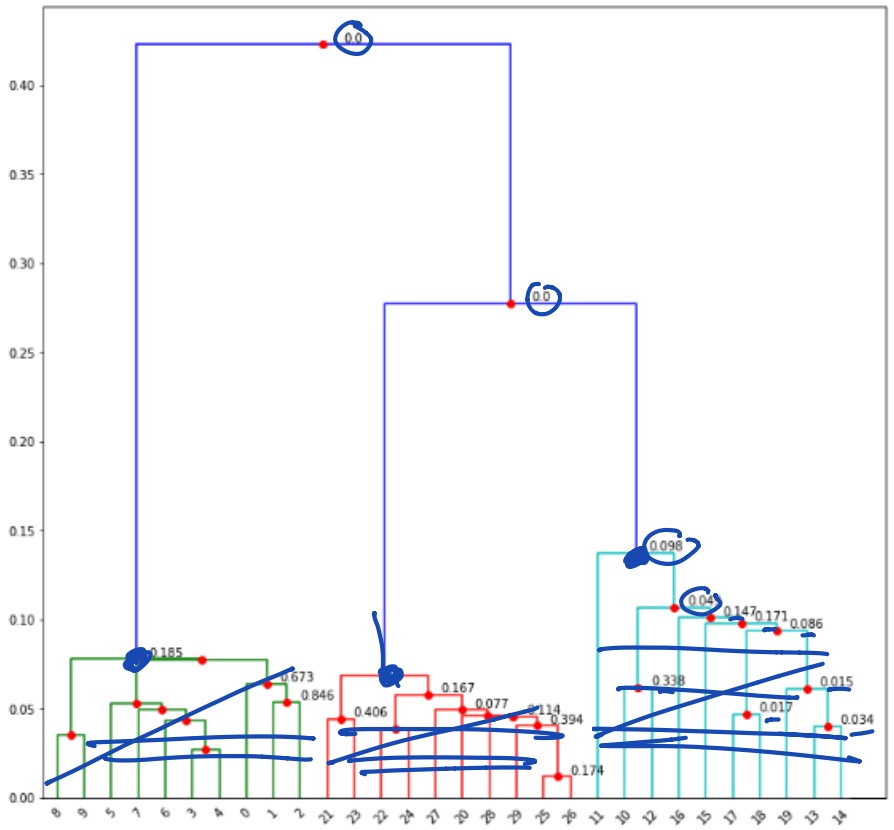


issues

- how many clusters
- are these clusters "significant"
- branch ordering







Centroid-Based Clustering

k-means clustering

For k clusters, find cluster members so that to minimize the distance to the cluster prototype (centroid).

NP-hard.

↳ k-means clustering

- choose k centroids, $u^{(i)}$
(randomly, heuristics)

- repeat

- assign objects to the nearest centroid

- move centroid to the center of the assigned objects

until no change in cluster assignments

goals

low complexity, k is low
speed, fast convergence

- clusters
wishes to be
"spherical"

problems

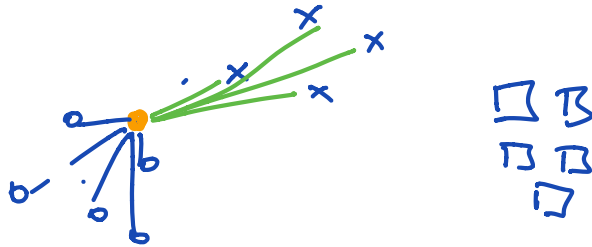
- depends on initialization
 k TRANSIT ||
- 2 to ...
- k defined

↙
how many clusters do we want?

Assessment of Cluster Quality

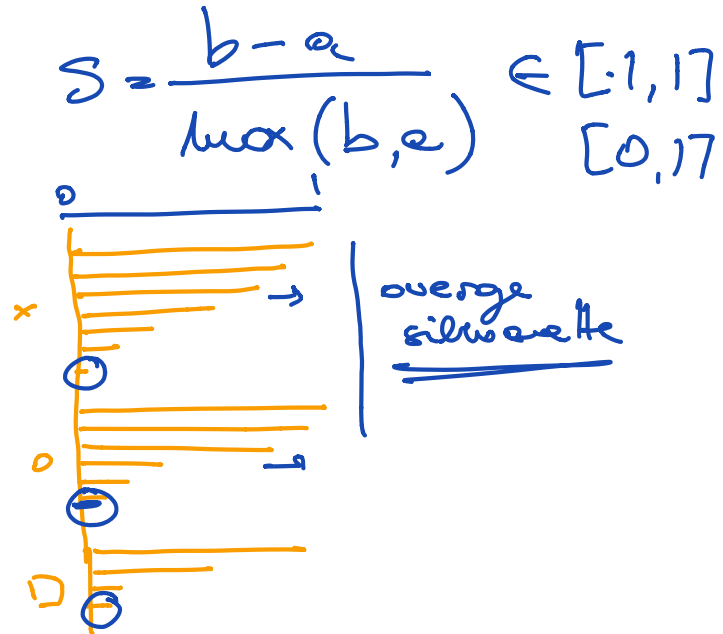
complexity \leftrightarrow cluster quality
↑ →

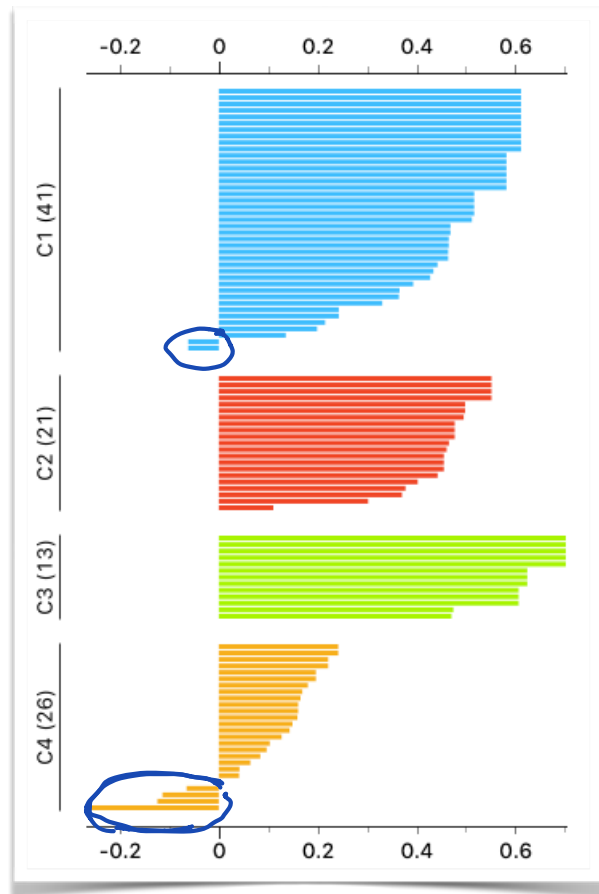
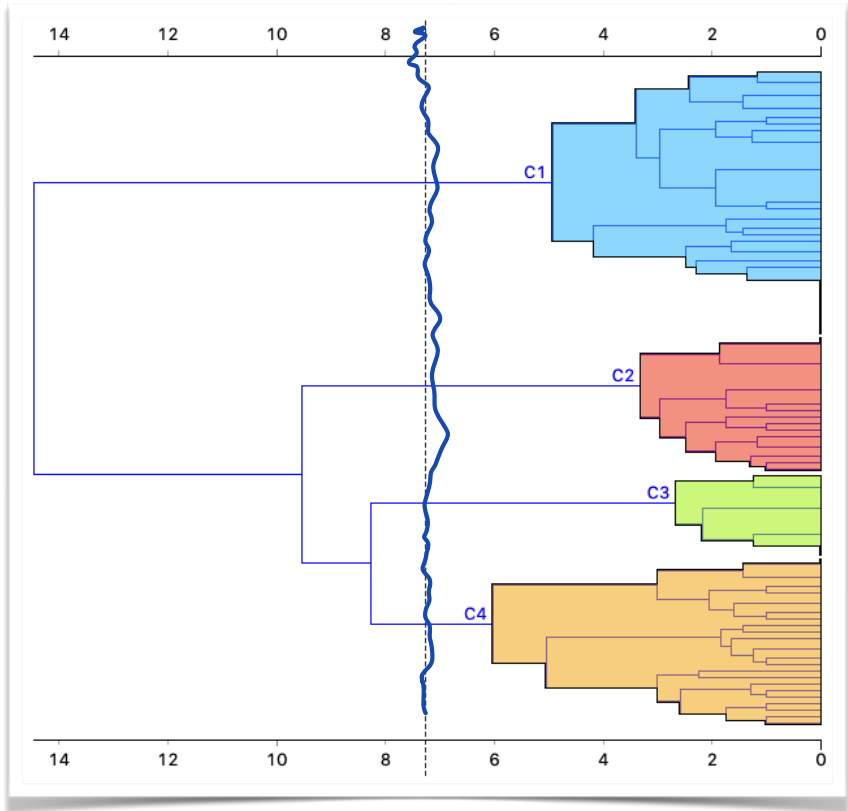
The silhouette

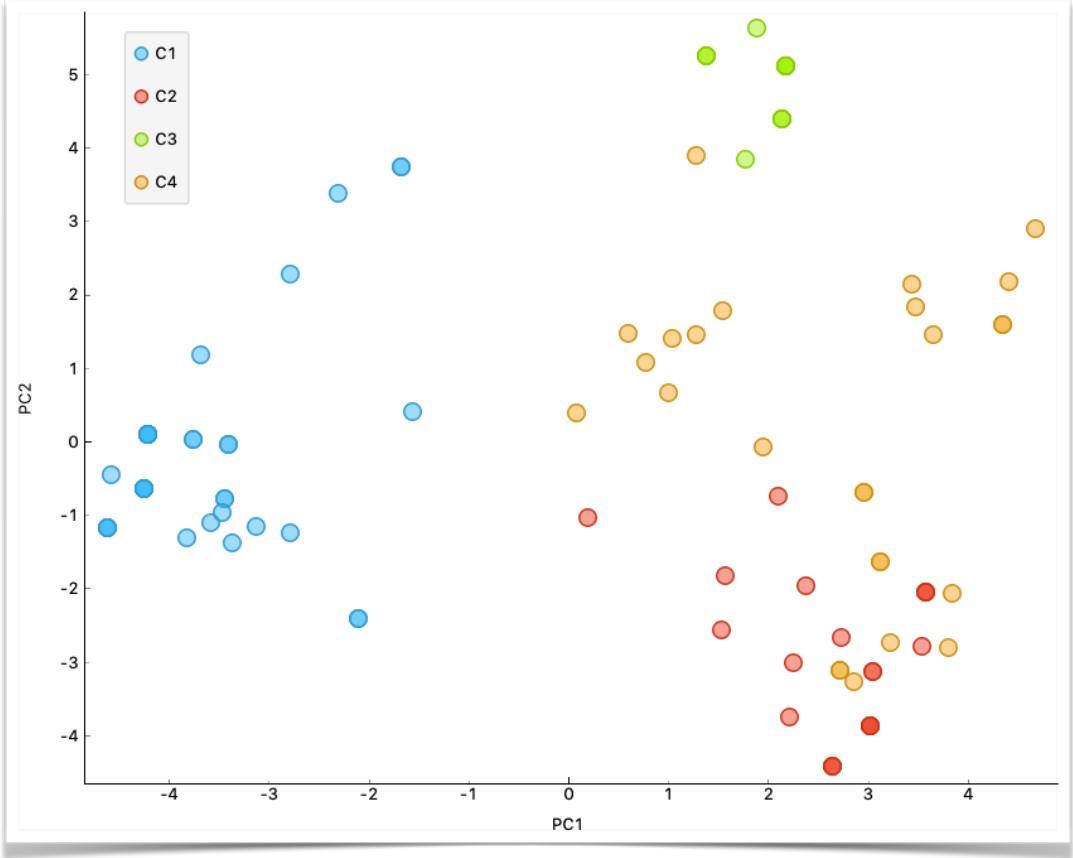


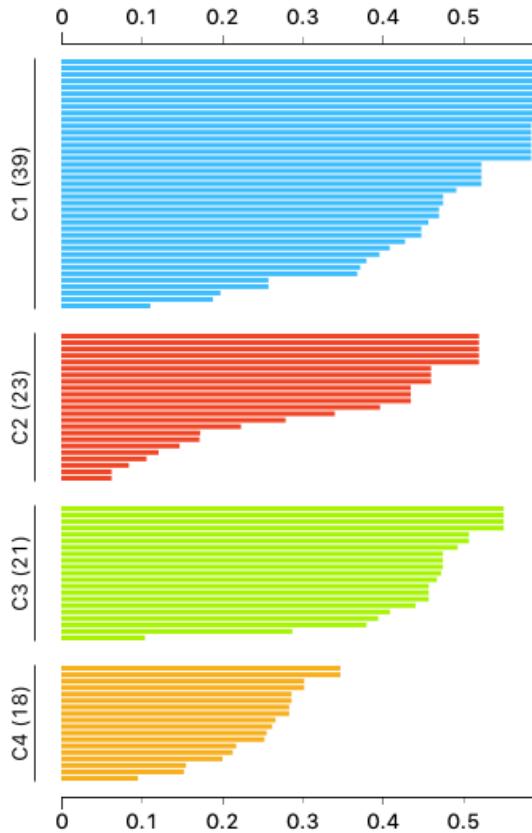
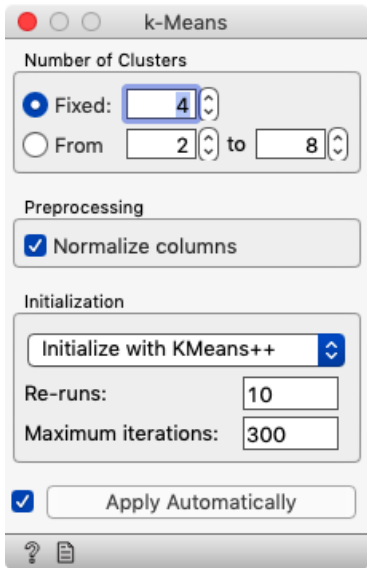
a : within cluster av. dist
 b : dist. to closest cluster

find k so that \bar{s} is maximized







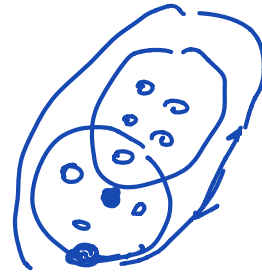


Density-Based Clustering

DBSCAN

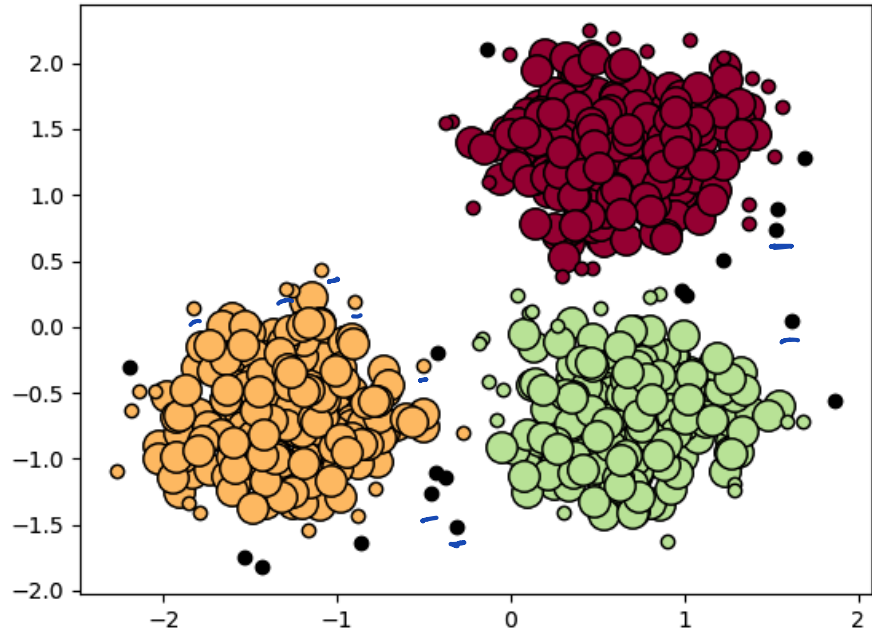
min Pts : minimal number of objects in the core
 ϵ : radius of the core

- find the cores
- merge overlapping cores

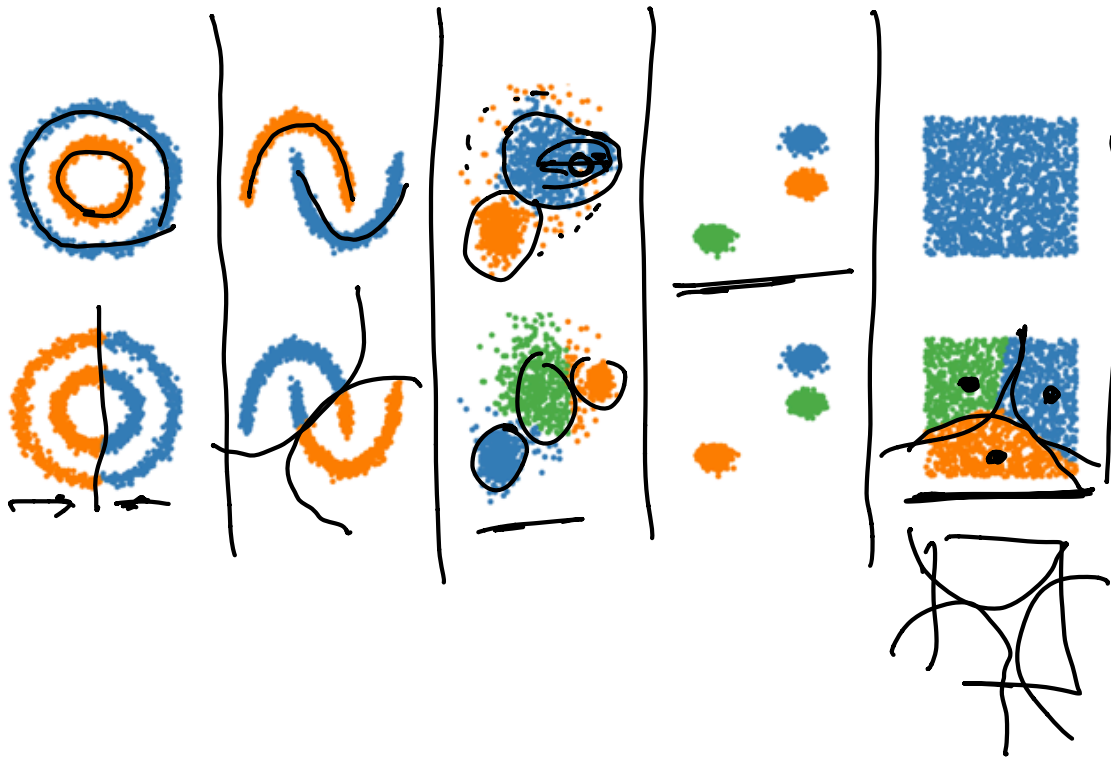


-
- not all points are clustered
 - some outliers of clusters
 - find dense regions

Estimated number of clusters: 3



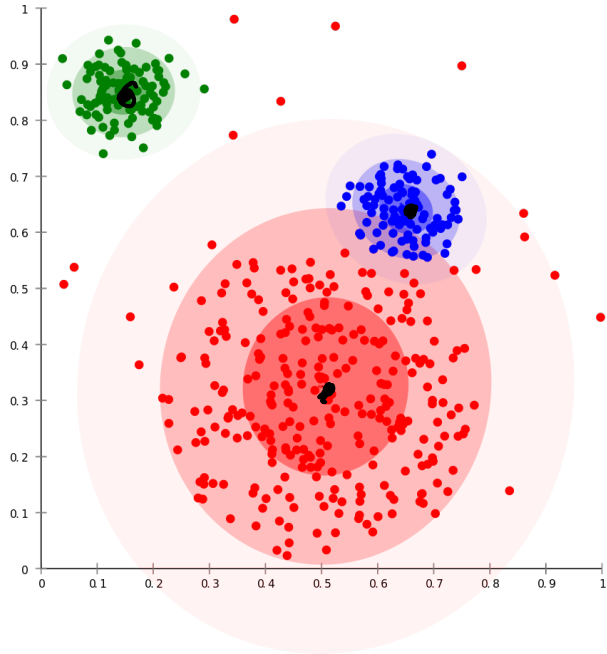
DBSCAN



k-means

Distribution-Based Clustering

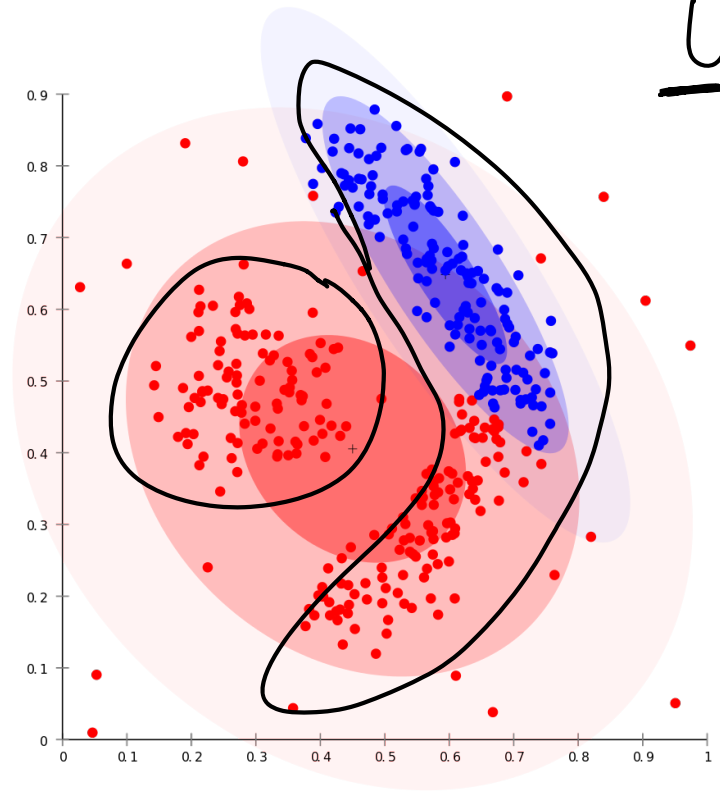
Probabilistic clustering



crisp clustering

fixed number of
Gaussian distributions
expectation maximization

УПДР



Interpretation

